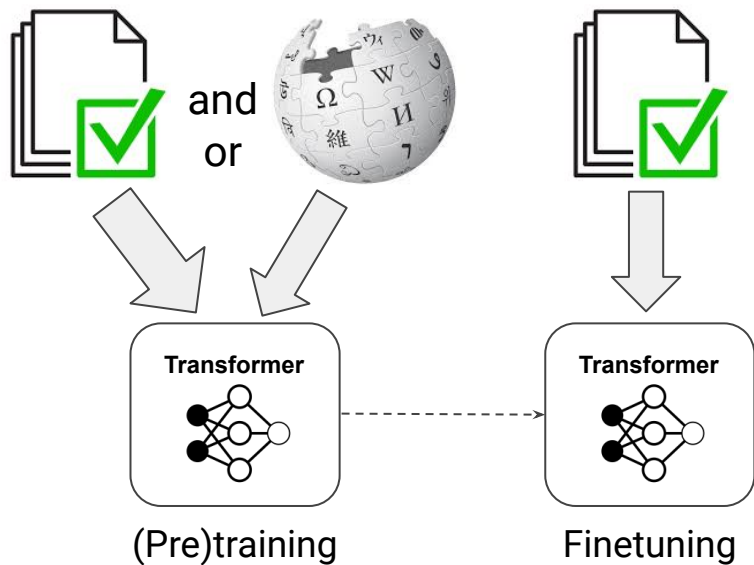


# Large-context and efficient models of language

Nikolaos Pappas

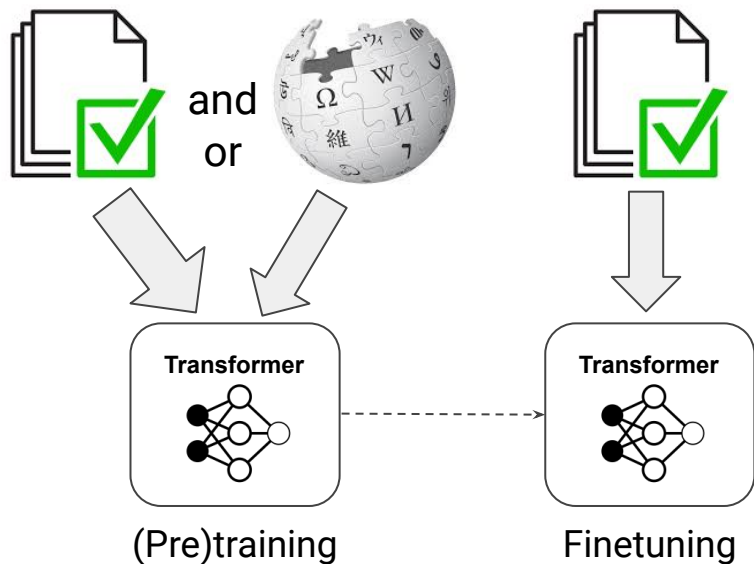


# Dominant paradigm



**Huge** impact in NLP tasks.

# Dominant paradigm



- ❌ Big data requirements
- ❌ Poor on rare or new words
- ❌ Computationally expensive

**Huge** impact in NLP tasks.

# My past research

Build models that learn from language **efficiently**

①

## Modeling documents

Representation learning  
[JAIR'17;EMNLP'18]

Multilingual transfer  
[IJNLP'17,TACL'19]

Structural comparisons  
[EMNLP'20]

②

## Promoting data efficiency

Deep word sharing  
[WMT'18;TACL'19;ICML'19]

Grounding to lexicons  
[EMNLP'20]

③

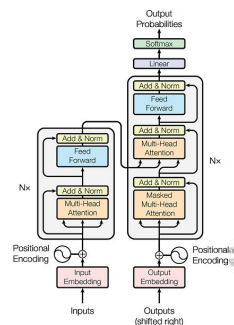
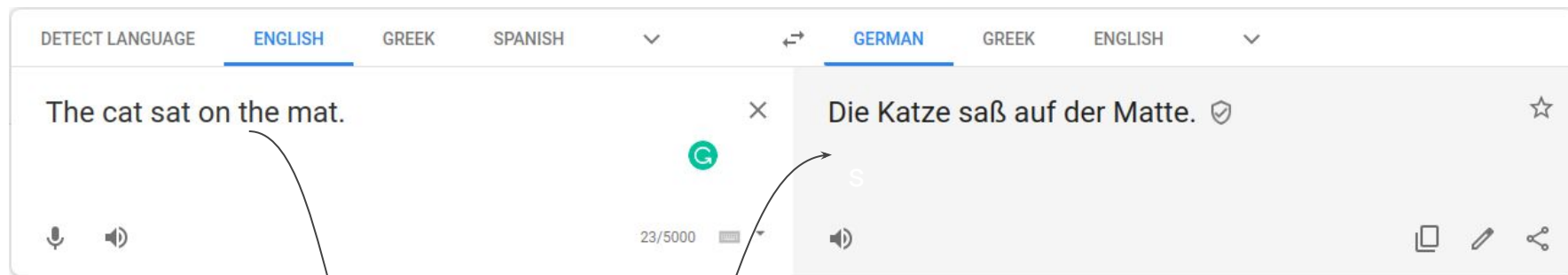
## Reducing cost

Model design  
[ICLR subm.]

Training objectives  
[EMNLP'20]

Scalable components  
[ICML'20, ICLR subm.]

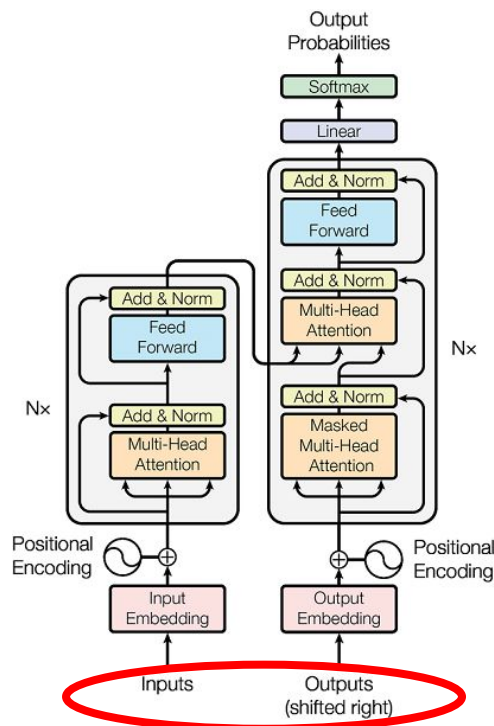
# Transformer origins



Transformer

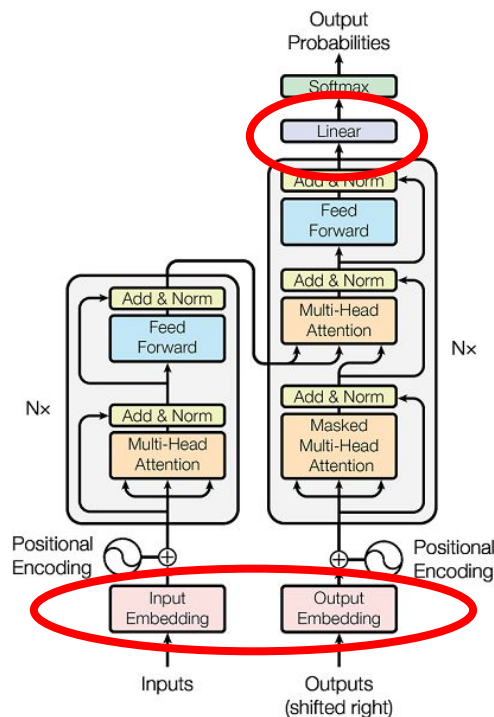
(Vaswani et al. 2017)

# Limitations: Narrow context



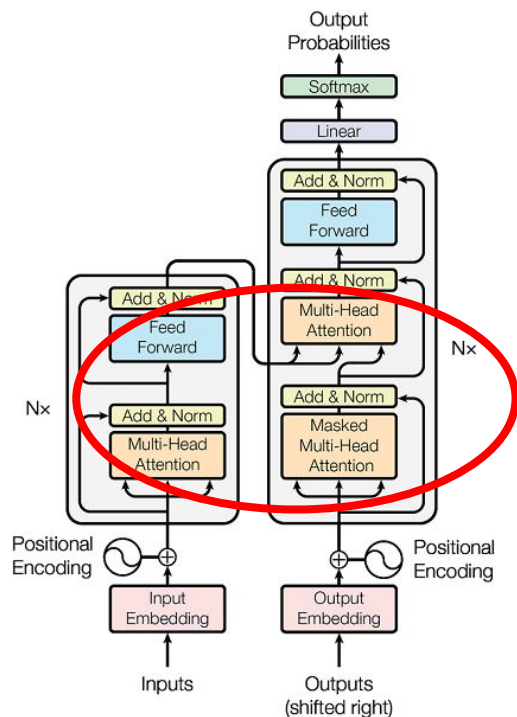
- Fixed and narrow context for prediction
- Suboptimal for document tasks

# Limitations: Rigid parameterization



- Dominates model size
- Performs poorly on rare types (data hungry)
- Requires ungraceful changes for adaptation

# Limitations: Quadratic complexity



- Does not scale to long text sequences
- Wastes memory for parallelization
- Slow for autoregressive inference



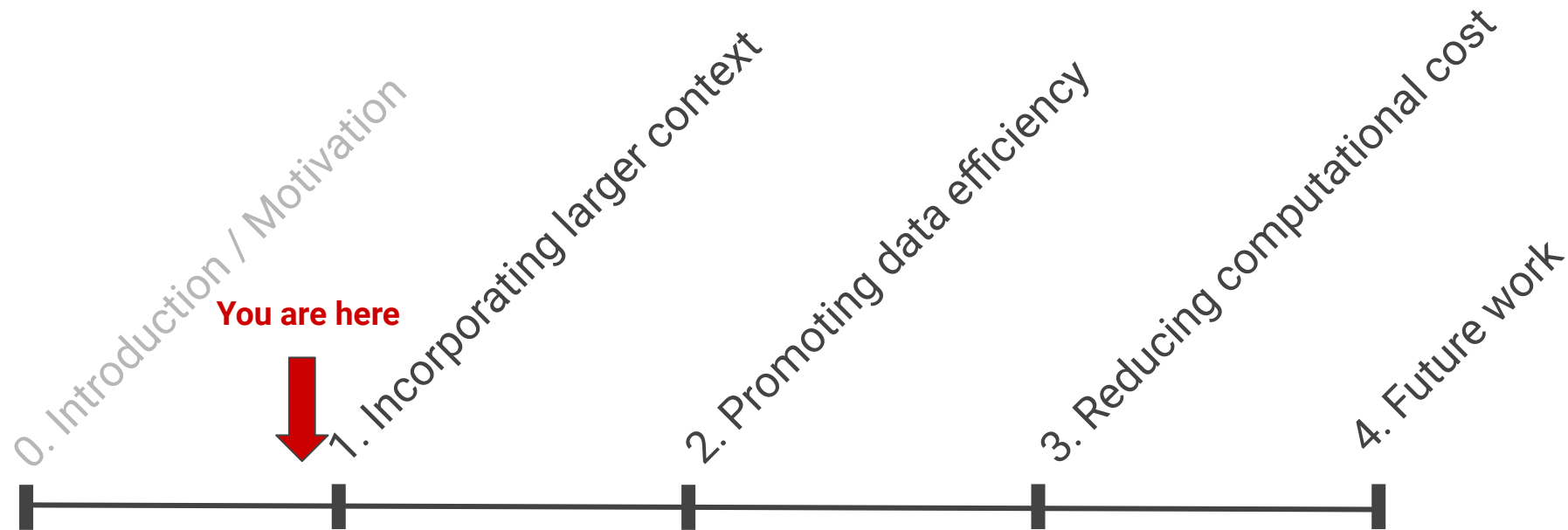
# Overview

Objective: show ways to address these challenges in neural MT

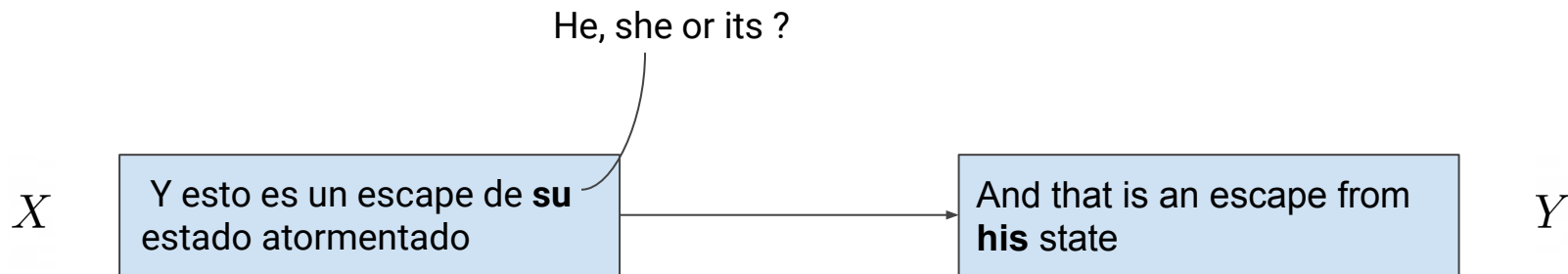


1. Dynamic hierarchical attention [[EMNLP 2018](#)]
2. Deep word sharing and grounding [[ICML 2019](#); [EMNLP 2020](#)]
3. Random feature attention [[ICLR subm.](#)]

# Overview



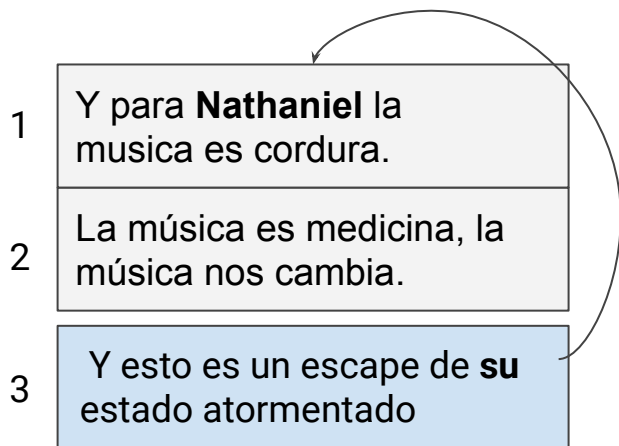
# Sentence-level translation



$$P(Y|X) = \prod_{t=1}^T p(y_t | y_{<t}, X)$$



# Previous efforts



- Concatenation (Tiedemann & Scherrer, 2017)
- Additional attention (Jean et al., 2017)
- Hierarchical context (Wang et al., 2017)
- Continuous cache (Tue et al., 2018)
- ... and many other recently (Voita et al., 2018; Lopes et al., 2020; Liu et al., 2020; Yu et al., 2020)

# Dynamic hierarchical context [EMNLP 2018]

- Exploit source and target document context
- Compute a dynamic document context for each token
- Increased interpretability in the attention maps

## Currently Translated Sentence

Src.: y esto es un escape de **su** estado atormentado .  
Ref.: and that is an escape from **his** tormented state .  
Base: and this is an escape from **its** < unk > state .  
Cache: and this is an escape from **their** state .  
HAN: and this is an escape from **his** < unk > state .

## Context from Previous Sentences

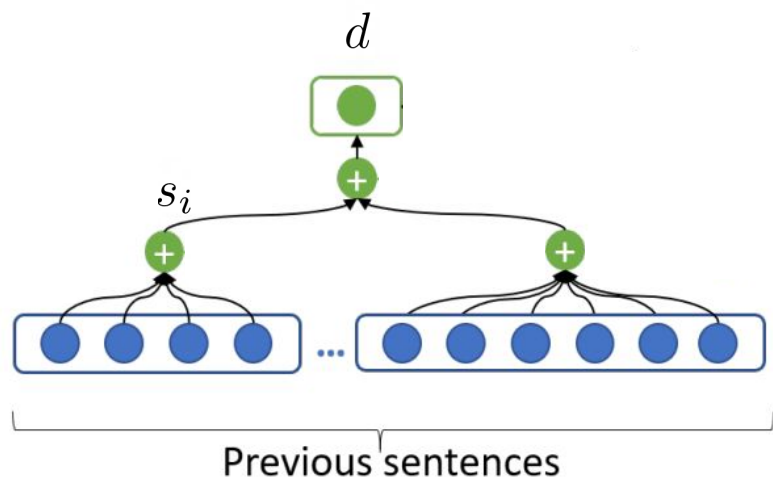
HAN decoder context with target. *Query: his* (En)

s<sup>t-3</sup> music is medicine . music changes us .  
s<sup>t-2</sup> and for Nathaniel , music is mine .  
s<sup>t-1</sup> because music allows him to take his thoughts and **his** delusions and **turn** through his imagination and his creativity actually .

HAN encoder context with source. *Query: su* (Es)

s<sup>t-3</sup> la música es medicina . la música nos cambia .  
s<sup>t-2</sup> y para **Nathaniel** la música es cordura .  
s<sup>t-1</sup> porque la música le permite tomar sus pensamientos y sus delirios y transformarlos a través de su imaginación y su creatividad en realidad .

# Hierarchical attention



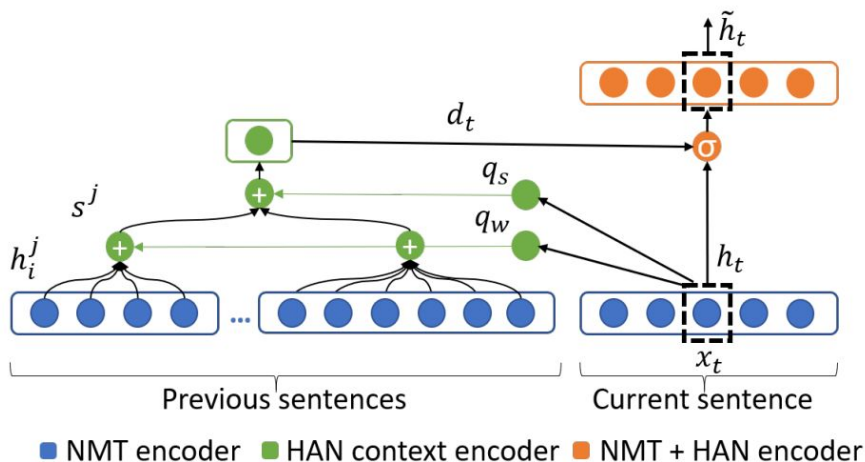
- Encoding with recurrent networks
- Pooling based on attention with a “learned context” per level

$$a_i = \frac{\exp(s_i^\top u_s)}{\sum_j \exp(s_j^\top u_s)}$$

$$d = \sum_i \alpha_i s_i$$

(Yang et al., 2017)

# Dynamic hierarchical attention [EMNLP 2018]



- Encoding with transformer  $s$
- Pooling based on multi-head attention conditioned on encoded tokens

$$q_s = f_s(h_t)$$

$$d_t = \text{FFN}(\text{MultiHead}_j(q_s, s^j))$$

- Context gating

$$\lambda_t = \sigma(W_h h_t + W_d d_t)$$

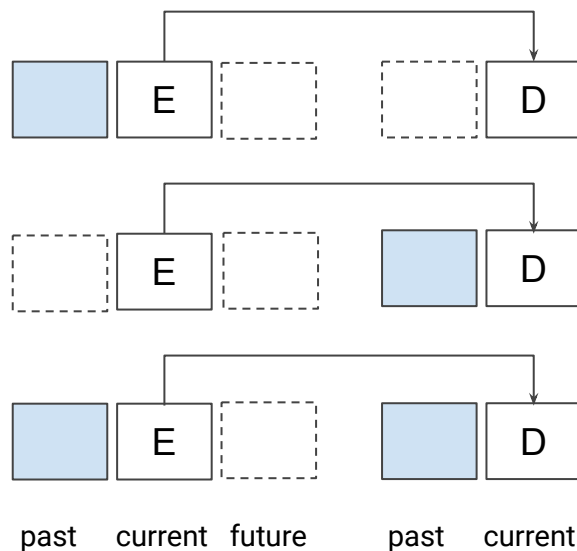
$$\tilde{h}_t = \lambda_t h_t + (1 - \lambda_t) d_t$$



# Document machine translation

- Document evaluation metrics
  - Noun/Pronoun accuracy
  - Lexical coherence: metric-based (LSA)
  - Lexical cohesion: repeated/total
- Datasets with document boundaries

TED Talks		Subtitles		News
Zh-En	Es-En	Zh-En	Es-En	Es-En
0.2M	0.2M	2.2M	4.0M	0.2M



# Sentence-level results

	TED Talks		Subtitles		News
	Zh-En	Es-En	Zh-En	Es-En	Es-En
NMT transformer	16.87	35.44	28.60	35.20	21.36
+ cache	17.32 ***	36.46 ***	28.86	35.49	22.36 ***
+ HAN encoder	17.61 *** ††	36.91 *** ††	29.35 * †	35.96 * †	22.36 ***
+ HAN decoder	17.39 ***	37.01 *** †††	29.21 * †	35.50	22.62 *** †††
+ HAN joint	<b>17.79</b> *** †††	<b>37.24</b> *** †††	<b>29.67</b> ** †	<b>36.23</b> ** ††	<b>22.76</b> *** †††

Higher is better

BLEU scores. Significance with respect to NMT \*, and to cache model †.

- Significant improvement on different size datasets (up to 4M)
- Target and source context are complementary

# Discourse-level results

	<b>Coherence</b>	<b>Lexical Cohesion</b>	<b>Pronouns</b>	<b>Nouns</b>
NMT transformer	28.42	47.98	62.84	52.50
+ HAN encoder	28.60	48.35	64.48	53.61
+ HAN decoder	28.78	48.51	64.04	53.55
+ HAN joint	28.82	48.61	64.32	54.19
Human reference	29.79	52.94	100.0	100.0

Higher is better

- Gains across the board especially for noun/pronoun translation
- Still a big gap between human reference and translations

# Takeaways [EMNLP 2018]

- Incorporating larger context with dynamic hierarchical attention
- Improves both sentence and discourse evaluation metrics
- Provides interpretability in the attention maps for each token

## Currently Translated Sentence

Src.: y esto es un escape de **su** estado atormentado .  
Ref.: and that is an escape from **his** tormented state .  
Base: and this is an escape from *its* < unk > state .  
Cache: and this is an escape from *their* state .  
HAN: and this is an escape from **his** < unk > state .

## Context from Previous Sentences

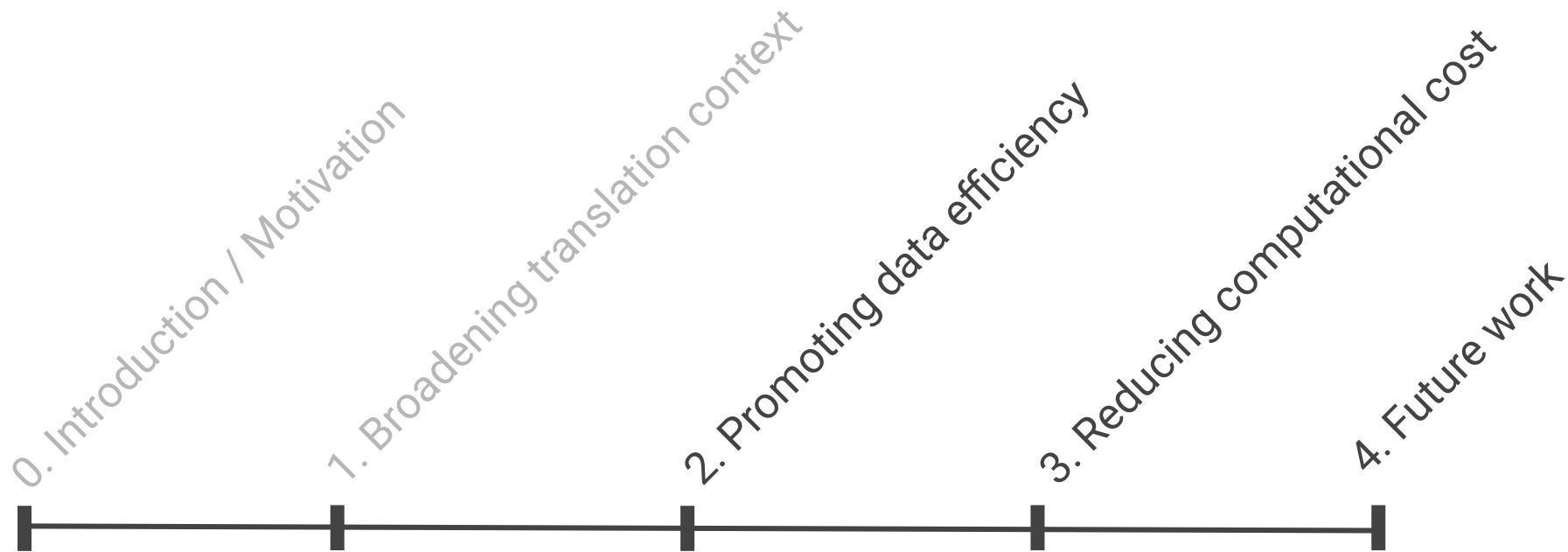
HAN decoder context with target. *Query: his* (En)

s<sup>t-3</sup> music is medicine . music changes us .  
s<sup>t-2</sup> and for Nathaniel , music is mine .  
s<sup>t-1</sup> because music allows him to take his thoughts and **his**  
delusions and **turn** through his imagination and his creat  
ivity actually .

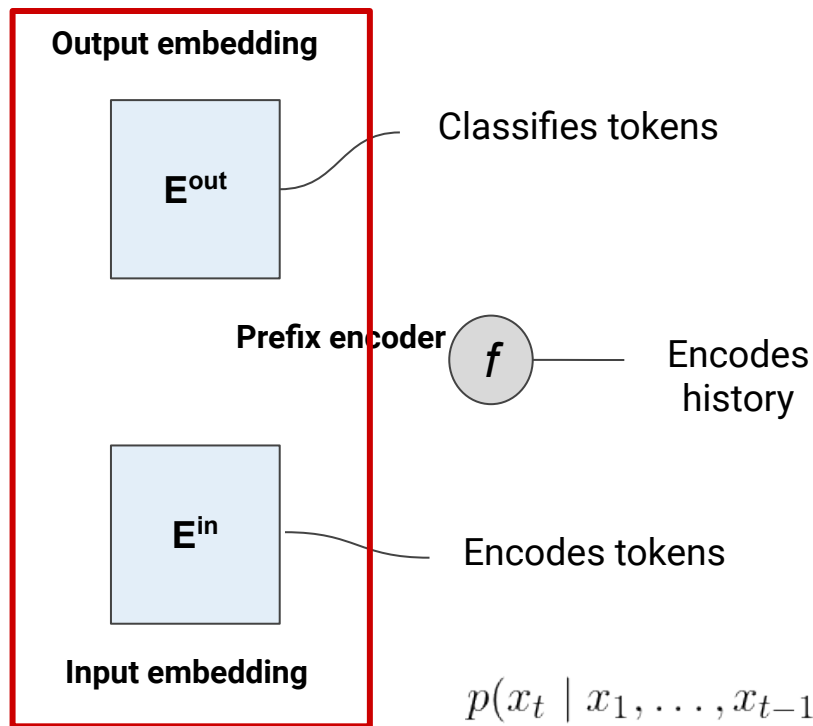
HAN encoder context with source. *Query: su* (Es)

s<sup>t-3</sup> la música es medicina . la música nos cambia .  
s<sup>t-2</sup> y para Nathaniel la música es cordura .  
s<sup>t-1</sup> porque la música le permite tomar sus pensamientos y  
sus delirios y transformarlos a través de su imaginació  
n y su creatividad en realidad .

# Overview



# Decoder language model

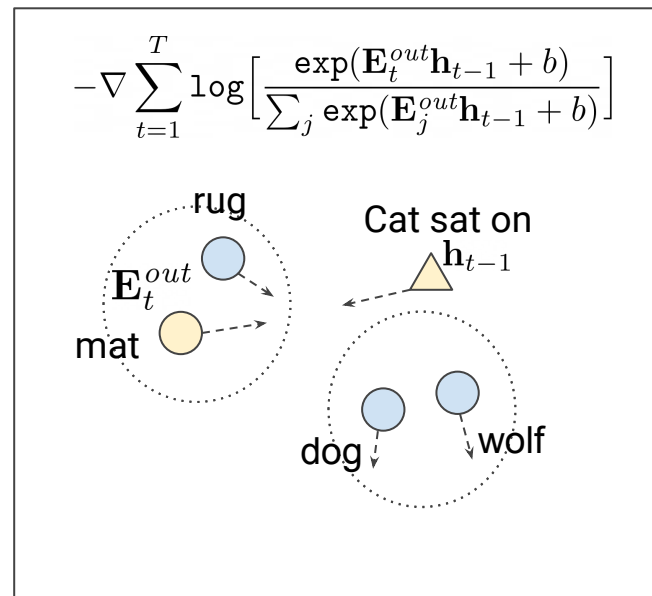


- Parameters are not shared across words
- Handle rare words poorly
- Cannot generalize well to new words or domains

$$p(x_t | x_1, \dots, x_{t-1}) \propto \exp(\mathbf{E}^{out} \mathbf{h}_{t-1} + \mathbf{b})$$

# Word sharing





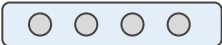
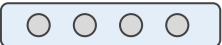


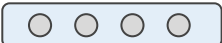



- Captures better the output space similarity
- Influences word neighbors during a training update



Net effect of training signal

(Pappas et al., 2018)

# Representing words

	Input		Output	Word sharing	Control power
Lookup table (Zaremba et al., 2014)		$\mathbf{E}^{out} \neq \mathbf{E}^{in}$			
Tied lookup table (Press and Wolf., 2017)		$\mathbf{E}^{out} = \mathbf{E}^{in}$			
Functional forms (Pappas et al., 2018; Gulurdova et al., 2018;) Bilinear Dual nonlinear		$\mathbf{E}^{out} = g(\mathbf{E}^{in})$			

- Increase power only via dim or rank which has the tendency to overfit in certain domains



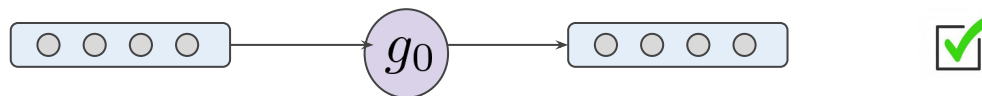
# Deep word sharing [ICML 2019]

Input

Output

Word sharing

$$\mathbf{E}_t^{out(1)} = g_0(\mathbf{E}_t^{in})$$



Functions  $g_j(\cdot)$  are simple nonlinear transformations

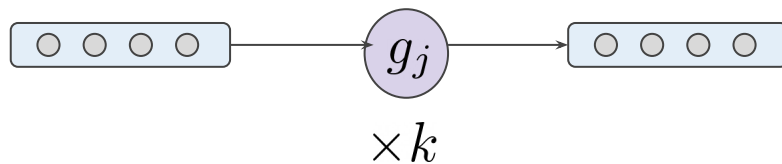
# Deep word sharing [ICML 2019]

Input

Output

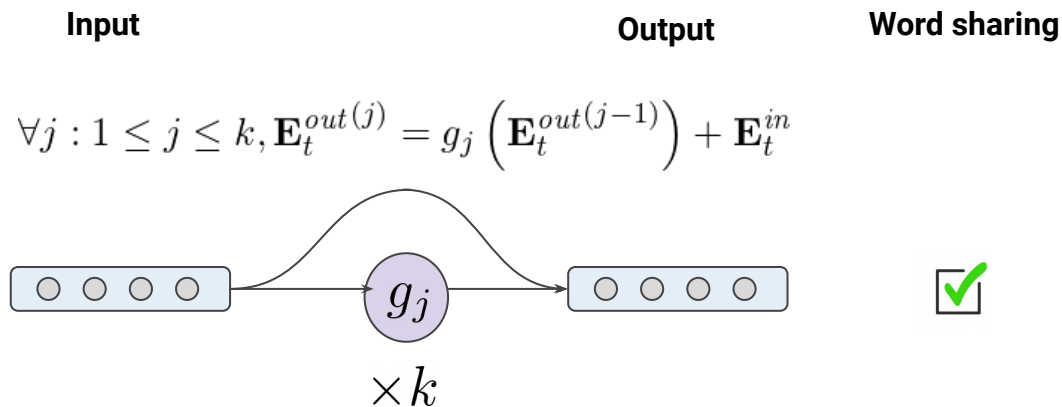
Word sharing

$$\forall j : 1 \leq j \leq k, \mathbf{E}_t^{out(j)} = g_j \left( \mathbf{E}_t^{out(j-1)} \right)$$



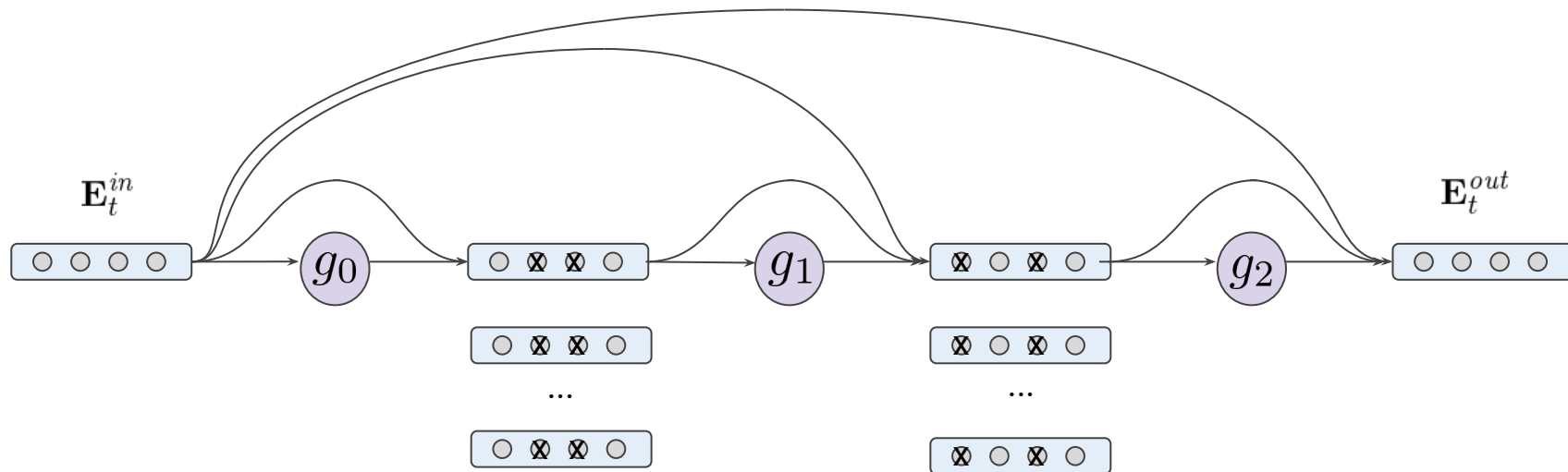
Functions  $g_j(\cdot)$  are simple nonlinear transformations

# Deep word sharing [ICML 2019]



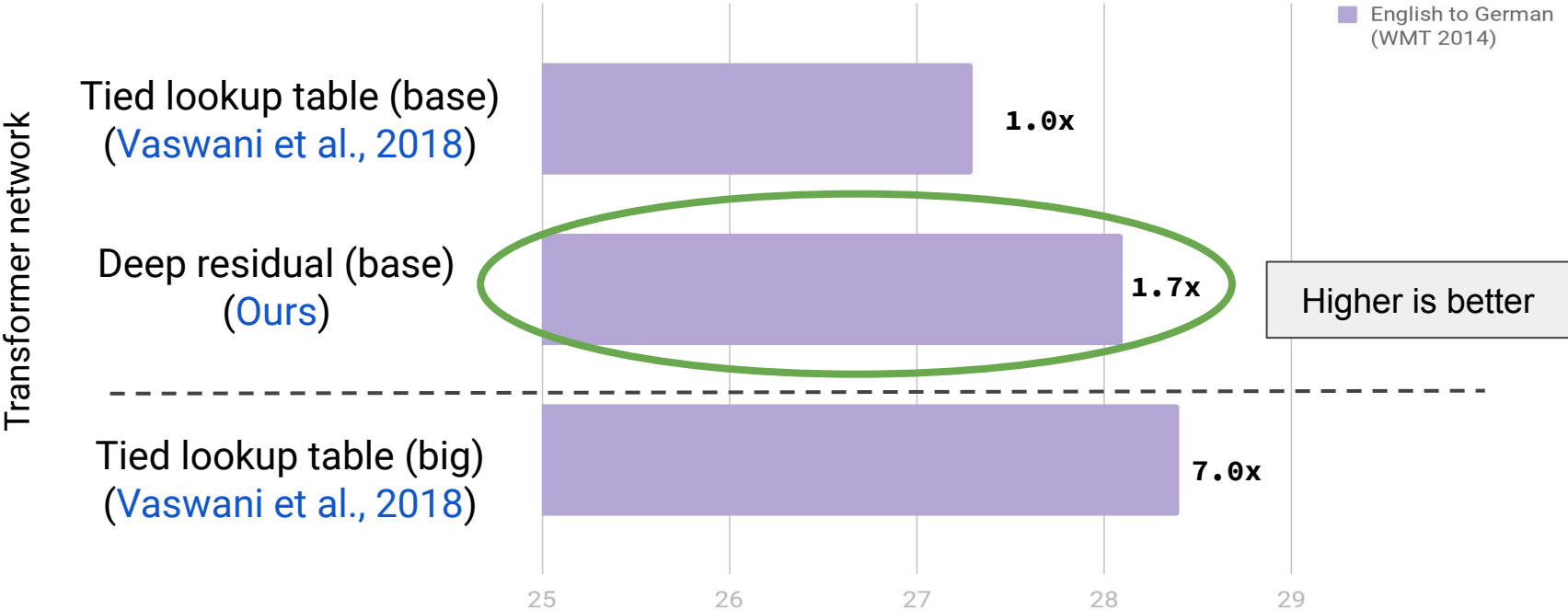
Functions  $g_j(\cdot)$  are simple nonlinear transformations

# Unfolded version with depth $k = 3$

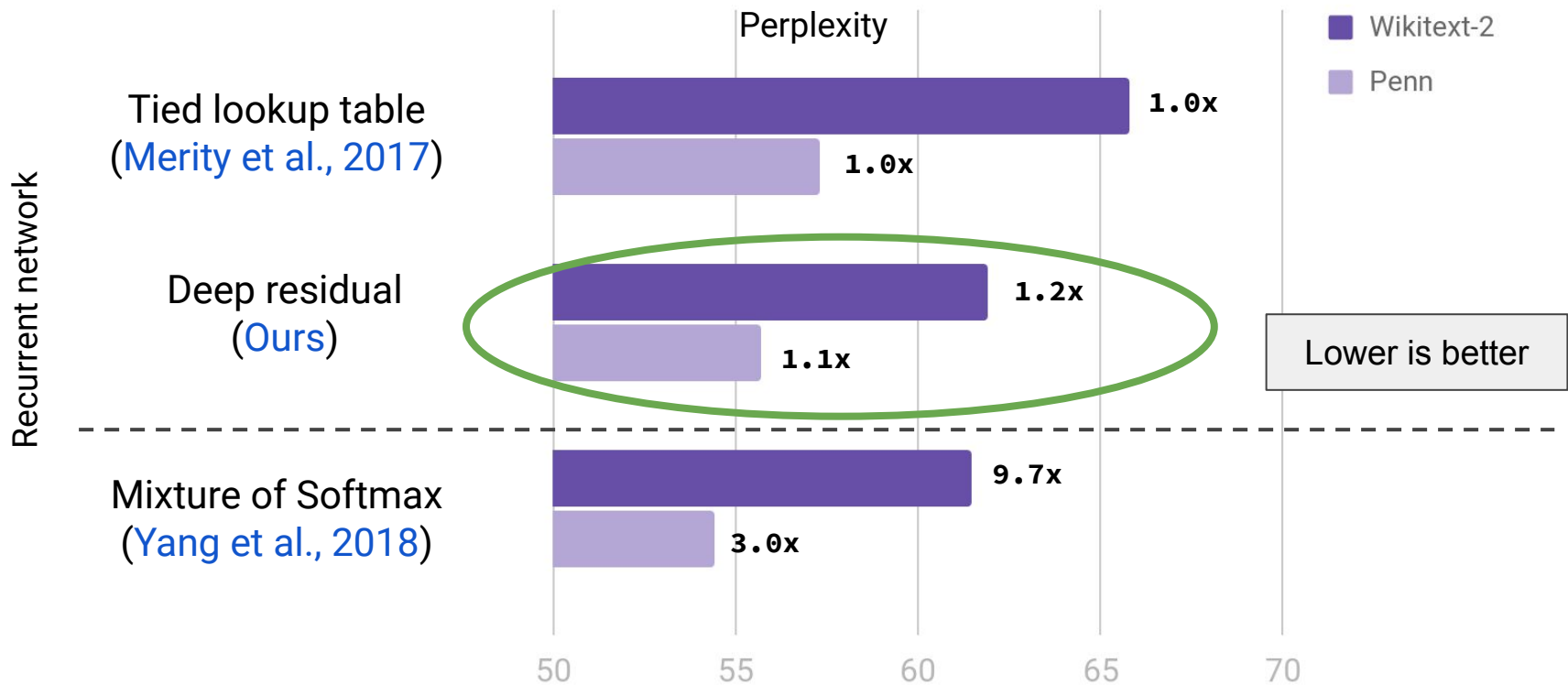


(Merity et al., 2017)

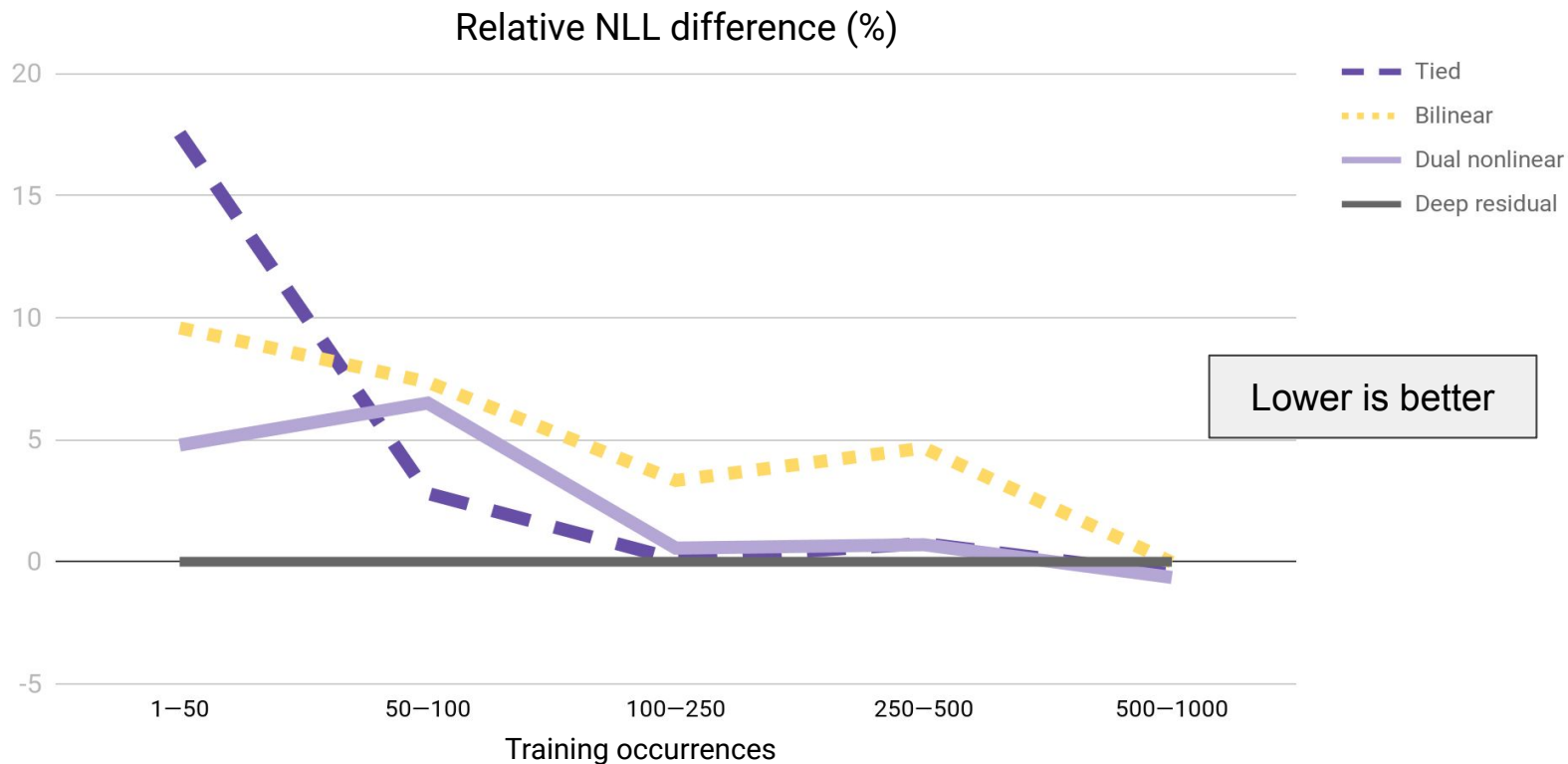
# Machine translation



# Language modeling



# Break down by frequency

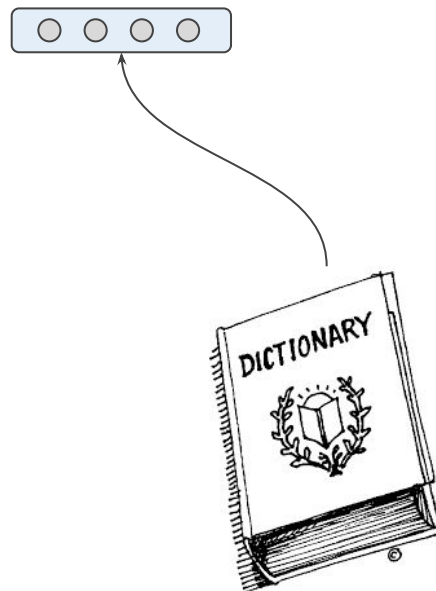


# Takeaways [ICML 2019]

- Deep word sharing improves speed-quality tradeoff
- Improvement is due to better modeling low-frequency words



Can we further gain by grounding to dictionaries and relaxing the vocabulary assumptions?





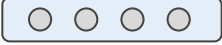
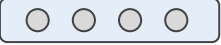


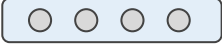
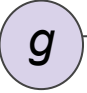
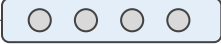


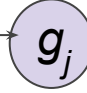






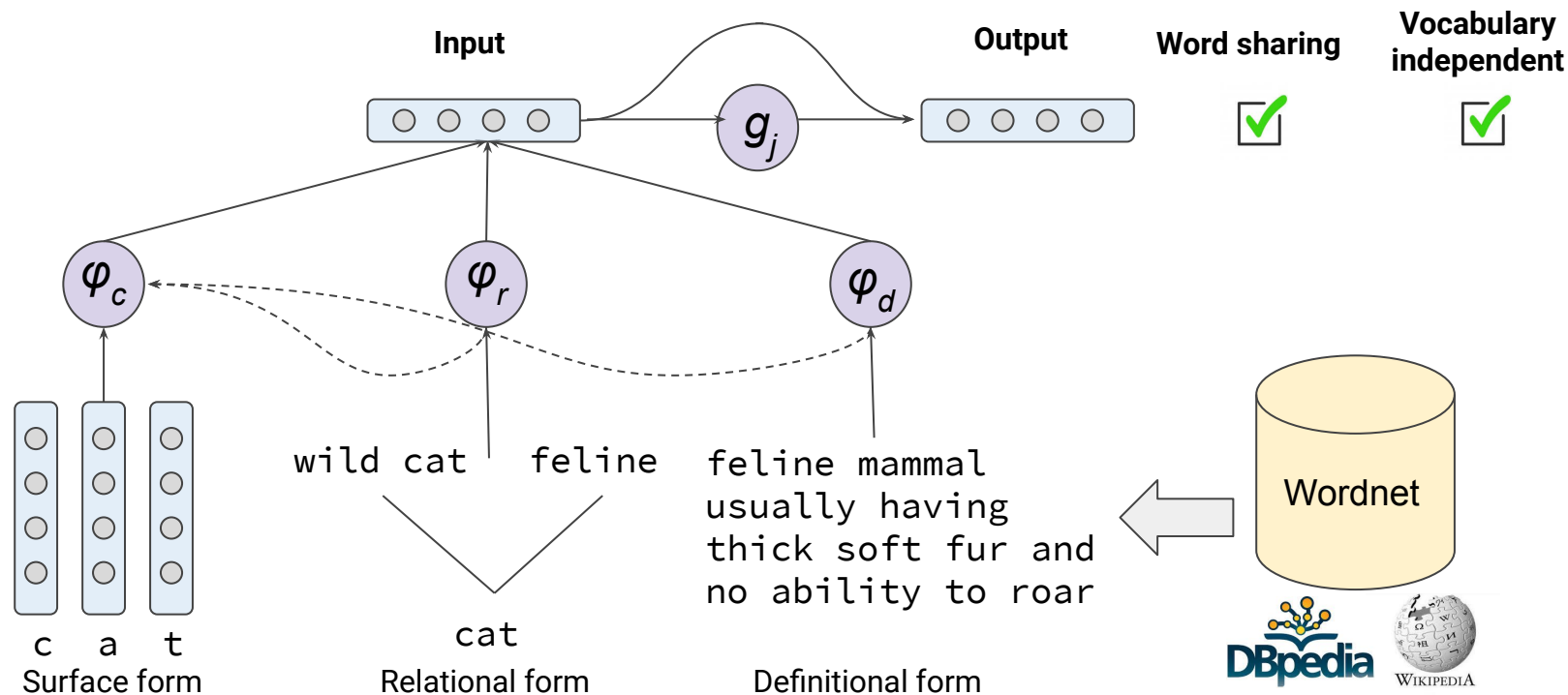
# Handling rare or new words

- Character-level models ([Cherry et al., 2018](#); [Al-Rfou et al., 2019](#))
  - ❌ Costly prefix encoders and training
- Data-driven vocabulary selection ([Sennrich et al., 2016](#); [Radford et al., 2018](#))
  - ❌ Linguistically simplistic
  - ❌ Rely on lookup tables
- Local neural cache ([Graves et al., 2017a,b](#))
  - ✅ Low-cost adaptation to rare/new words

# Related work: Word sharing

	Input		Output	Word sharing	Vocabulary independent
Lookup tables (Zaremba et al., 2014; Press & Wolf, 2017)		$\mathbf{E}^{out} \neq \mathbf{E}^{in}$			
		$\mathbf{E}^{out} = \mathbf{E}^{in}$			
		$\mathbf{E}^{out} = g(\mathbf{E}^{in})$			
Compositional / Functional forms (Jozefowicz et al., 2016; Baevski & Auli, 2019; Pappas & Henderson, 2019)					
		$\forall j : 1 \leq j \leq k, \mathbf{E}_t^{out(j)} = g_j(\mathbf{E}_t^{out(j-1)}) + \mathbf{E}_t^{in}$			

# Grounded compositional outputs [EMNLP 2020]



# Adapting to any vocabulary [EMNLP 2020]

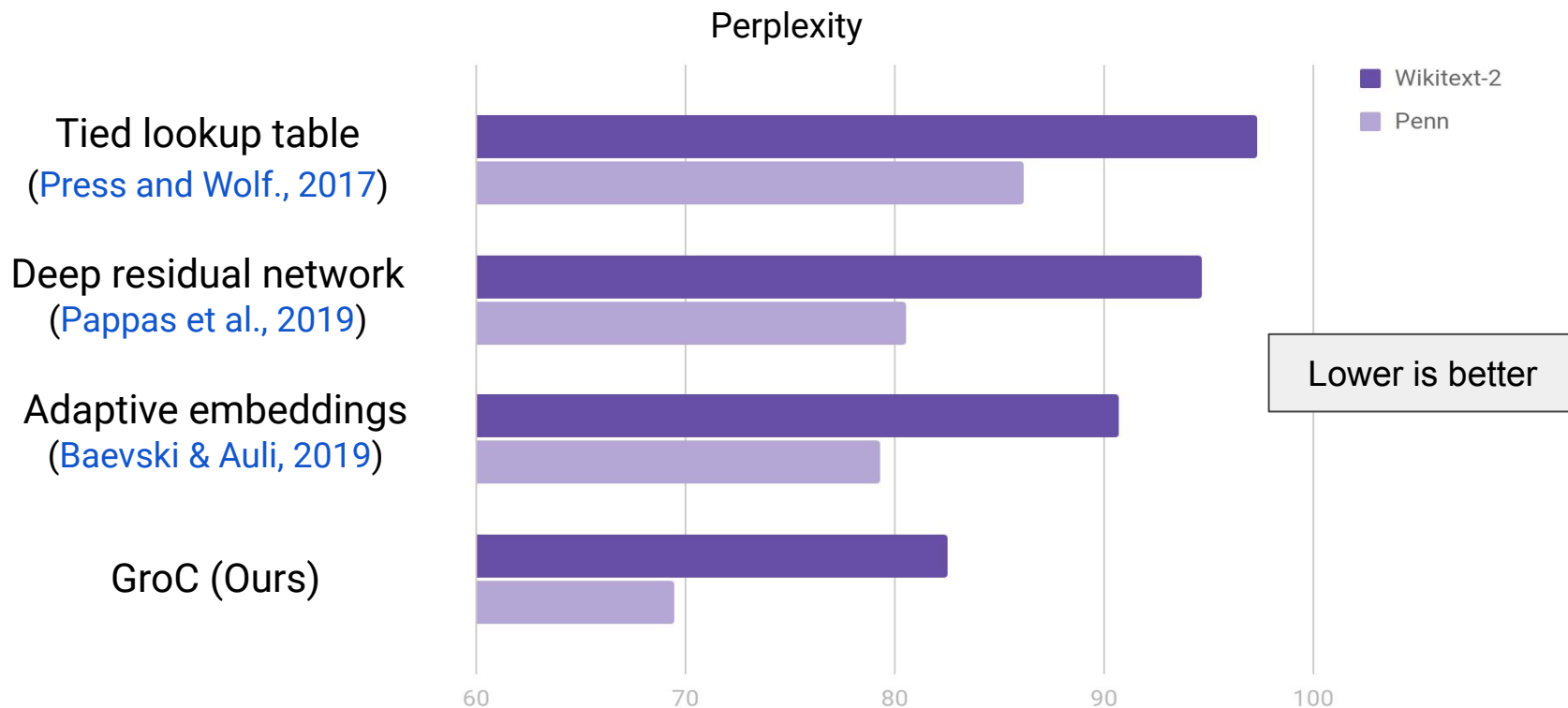
- We first represent the vocabulary with GroC
- Then we estimate the bias for each word  $u$

$$\mathbf{E}^{out} = \text{GroC}(\mathcal{V}^*)$$

$$b_v = \sigma(\mathbf{w} \cdot \mathbf{e}_v^{out} + a)$$

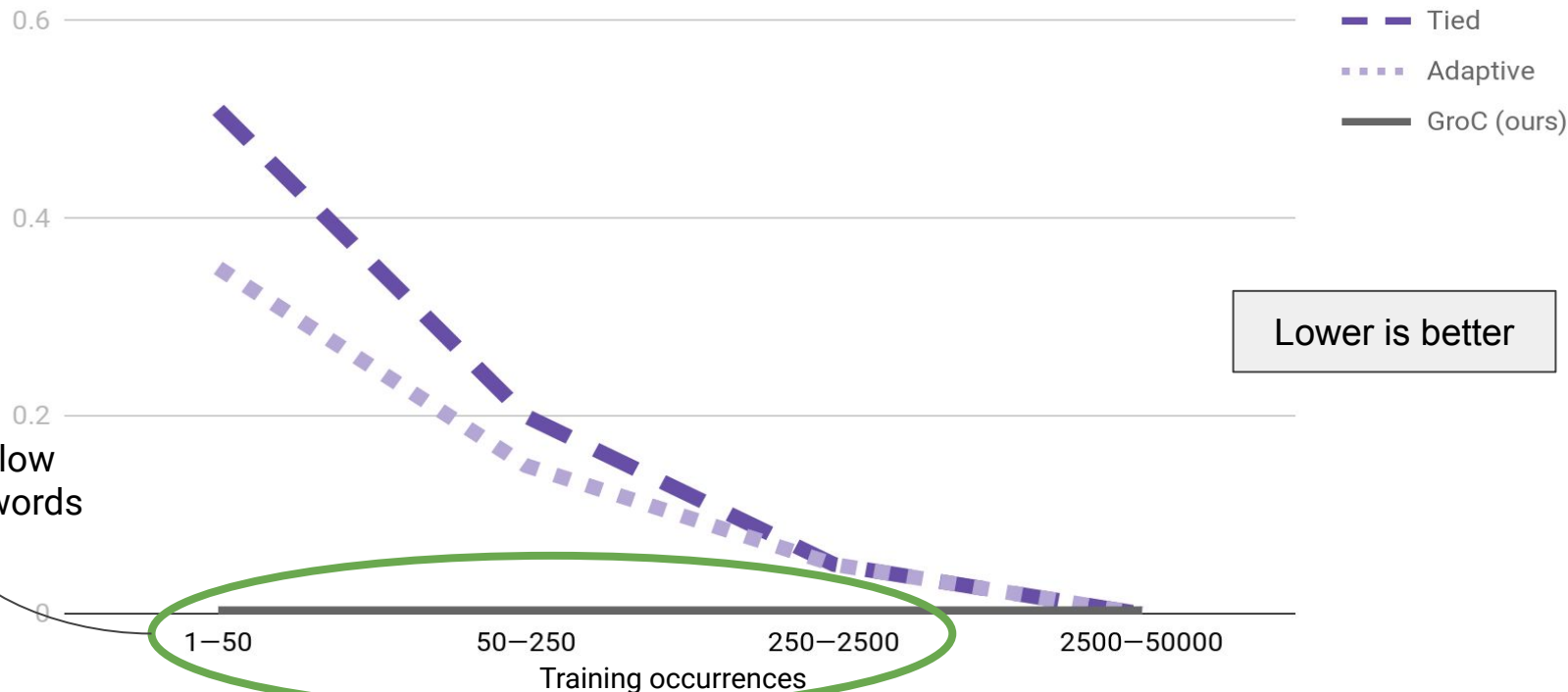

$$p(X_t = x_t \mid \mathbf{h}_{t-1}) \propto \exp(\mathbf{E}^{out} \mathbf{h}_{t-1} + \mathbf{b})$$

# Conventional language modeling

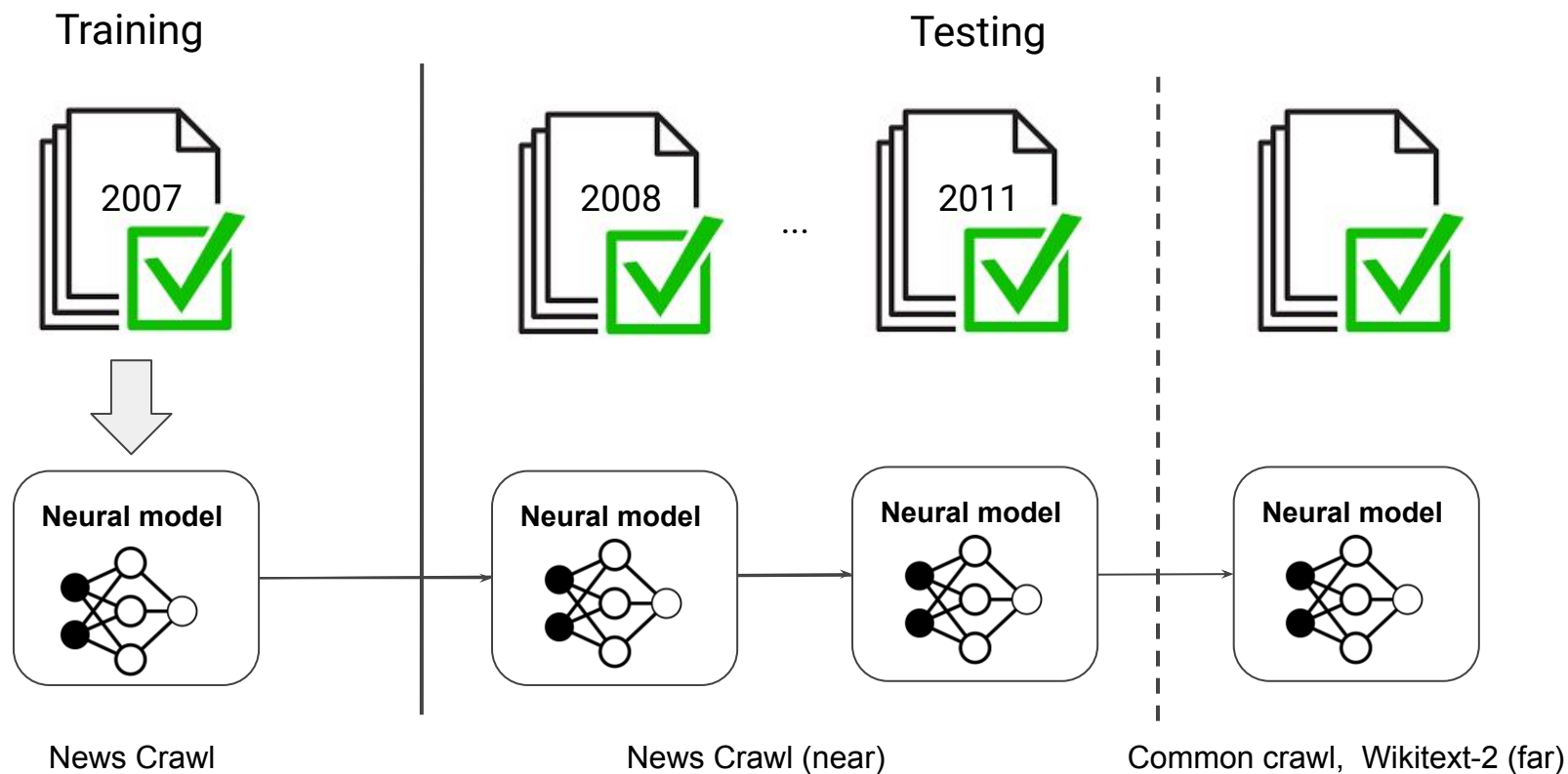


# Break down by frequency

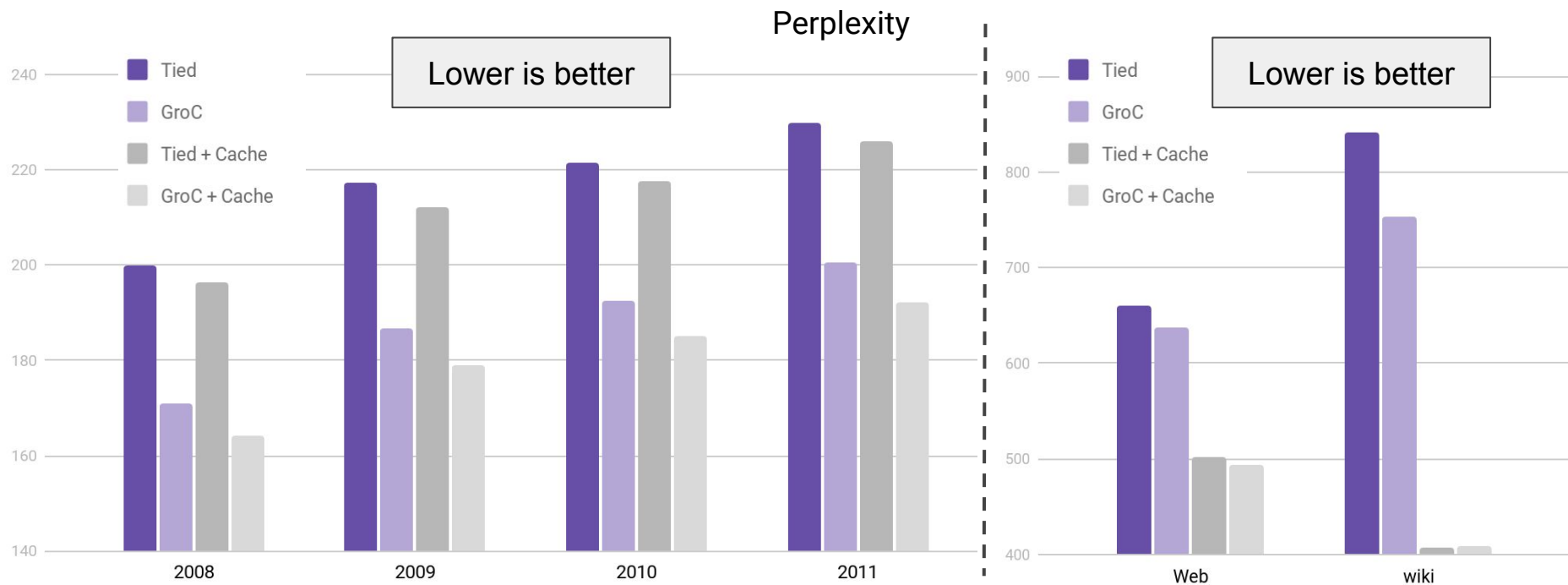
Median NLL difference



# Cross-domain modeling: Zero resources



# Cross-domain modeling: Zero resources

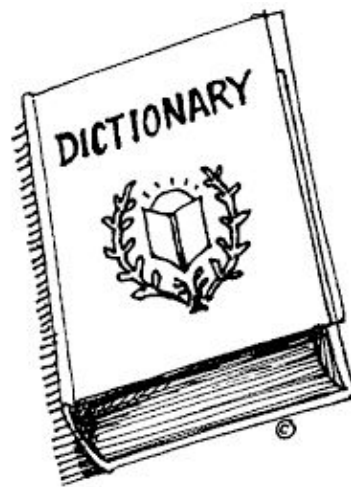




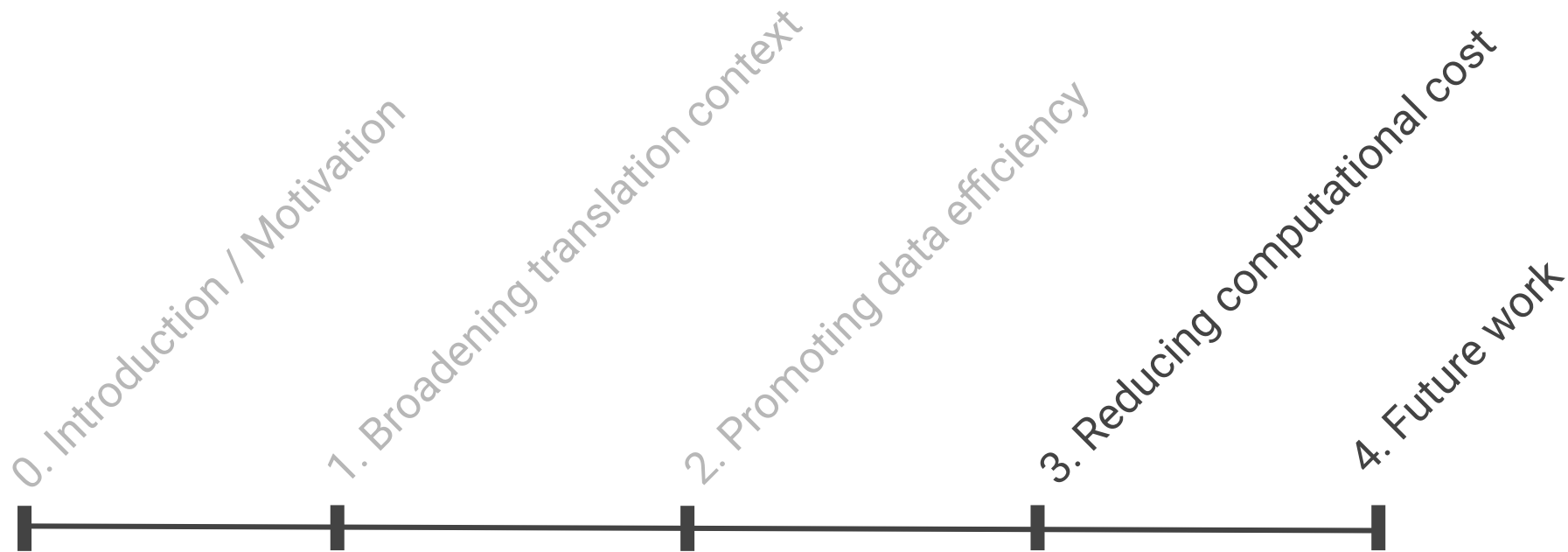
# Takeaways [EMNLP 2020]

## Grounded compositional word sharing

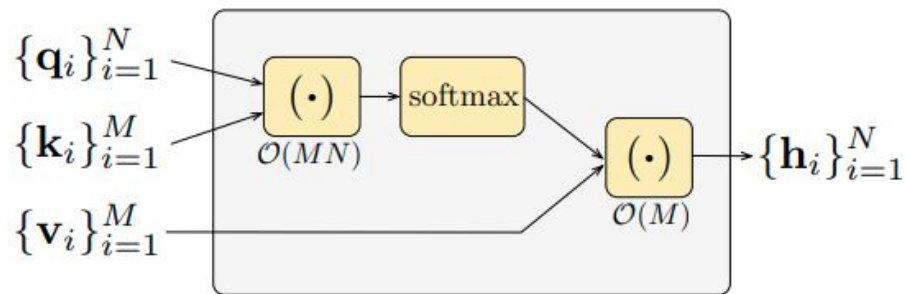
- Creates a compact representation of any vocabulary
- Achieves low perplexity on rare or new words
- Generalizes well to previously unseen domains



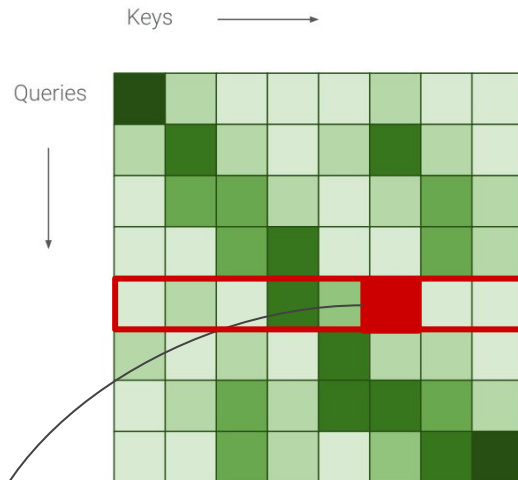
# Overview



# Softmax attention

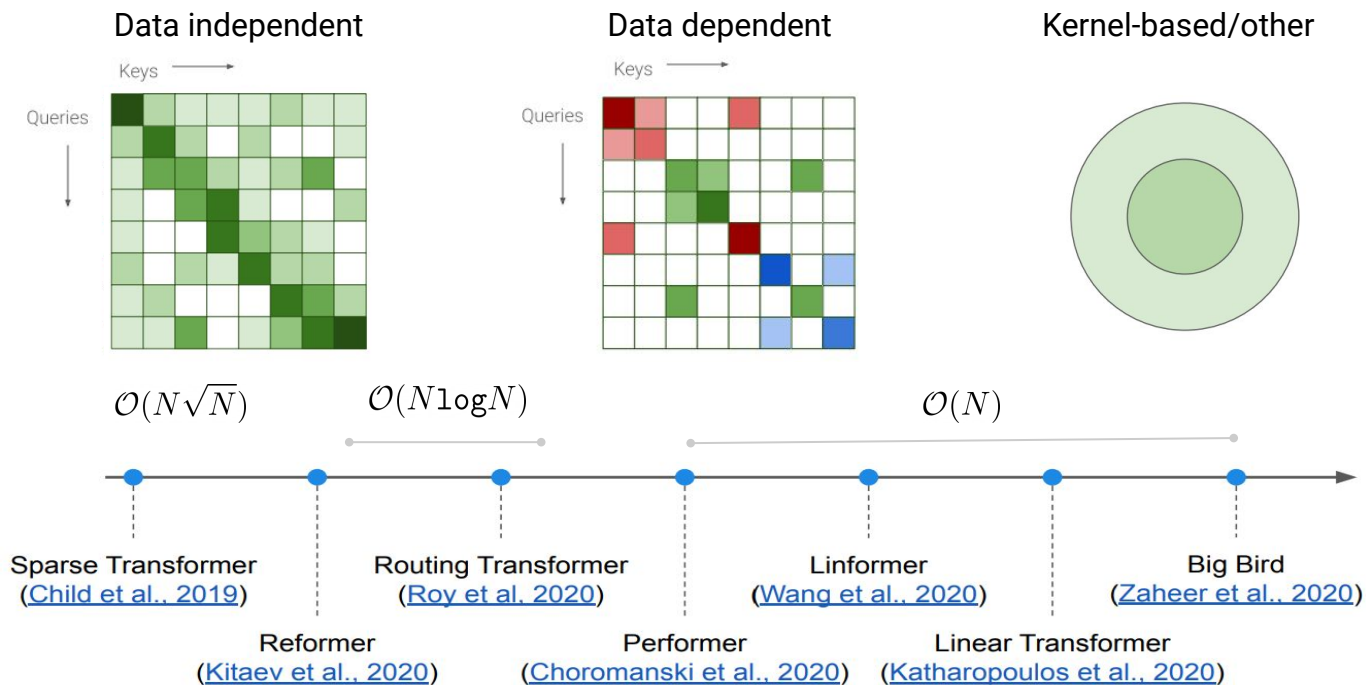


(a) Softmax attention.



$$\text{attn}(\mathbf{q}_t, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) = \sum_i \frac{\exp(\mathbf{q}_t \cdot \mathbf{k}_i / \tau)}{\sum_j \exp(\mathbf{q}_t \cdot \mathbf{k}_j / \tau)} \mathbf{v}_i^\top$$

# Recent progress



\* Ilharco et al., High-performance NLP tutorial, EMNLP 2020

# Results so far

Model / Paper	Complexity	Decode
Memory Compressed <sup>†</sup> (Liu et al., 2018)	$\mathcal{O}(n_c^2)$	✓
Image Transformer <sup>†</sup> (Parmar et al., 2018)	$\mathcal{O}(n.m)$	✓
Set Transformer <sup>†</sup> (Lee et al., 2019)	$\mathcal{O}(nk)$	✗
Transformer-XL <sup>†</sup> (Dai et al., 2019)	$\mathcal{O}(n^2)$	✓
Sparse Transformer (Child et al., 2019)	$\mathcal{O}(n\sqrt{n})$	✓
Reformer <sup>†</sup> (Kitaev et al., 2020)	$\mathcal{O}(n \log n)$	✓
Routing Transformer (Roy et al., 2020)	$\mathcal{O}(n \log n)$	✓
Axial Transformer (Ho et al., 2019)	$\mathcal{O}(n\sqrt{n})$	✓
Compressive Transformer <sup>†</sup> (Rae et al., 2020)	$\mathcal{O}(n^2)$	✓
Sinkhorn Transformer <sup>†</sup> (Tay et al., 2020b)	$\mathcal{O}(b^2)$	✓
Longformer (Beltagy et al., 2020)	$\mathcal{O}(n(k+m))$	✓
ETC (Ainslie et al., 2020)	$\mathcal{O}(n_g^2 + nn_g)$	✗
Synthesizer (Tay et al., 2020a)	$\mathcal{O}(n^2)$	✓
Performer (Choromanski et al., 2020)	$\mathcal{O}(n)$	✓
Linformer (Wang et al., 2020b)	$\mathcal{O}(n)$	✗
Linear Transformers <sup>†</sup> (Katharopoulos et al., 2020)	$\mathcal{O}(n)$	✓
Big Bird (Zaheer et al., 2020)	$\mathcal{O}(n)$	✗

- Transformers can be more memory and compute efficient
- Benefits mostly when they are trained on longer sequences
- Evaluation is often tricky

Higher accuracy on long sequence tasks

Lower perplexity on LM with longer context

Faster inference on CIFAR10

(Tay et al., 2020)

# Random feature attention (in a nutshell) [ICLR subm.]

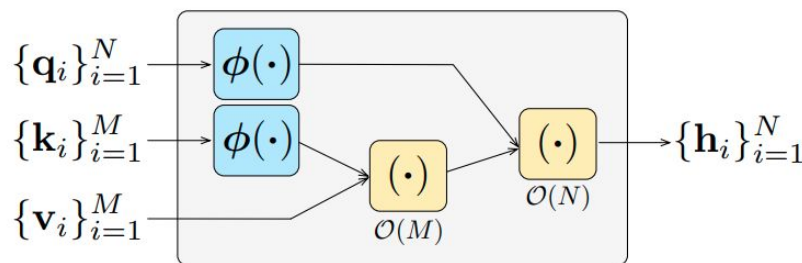
## ✓ Unbiased approximation of softmax attention

- 2X faster on MT decoding
- 17X faster on LM decoding
- 5X faster on long text classification

## ✓ Realistic speed/quality estimates

- Moderate and long sequence tasks
- Measurements with fixed batch size

## ✓ New insights on how to improve attention



$$\text{RFA}(\mathbf{q}_t, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) = \frac{\phi(\mathbf{q}_t)^\top \sum_i \phi(\mathbf{k}_i) \otimes \mathbf{v}_i}{\phi(\mathbf{q}_t) \cdot \sum_j \phi(\mathbf{k}_j)}$$

# Random Fourier features

- Can approximate a desired shift-invariant kernel e.g. Gaussian or Arccos
- Let  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$  be a nonlinear transformation and  $w_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ .

$$\phi(\mathbf{x}) = \sqrt{1/D} \left[ \sin(\mathbf{w}_1 \cdot \mathbf{x}), \dots, \sin(\mathbf{w}_D \cdot \mathbf{x}), \cos(\mathbf{w}_1 \cdot \mathbf{x}), \dots, \cos(\mathbf{w}_D \cdot \mathbf{x}) \right]^\top$$

then it provides an unbiased approximation of Gaussian kernel

$$\exp(-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2) \approx \phi(\mathbf{x}) \cdot \phi(\mathbf{y}) \quad (1)$$

(Rahimi & Recht, 2008)

# Random feature attention [ICLR subm.]

- Exponential becomes Gaussian if  $x, y$  are normalized (Rawat et al., 2019)

$$\exp(x \cdot y / \sigma^2) = \exp(1 / \sigma^2) \exp(-\|x - y\|^2 / 2\sigma^2)$$

- Therefore RFA is derived as follows

$$\begin{aligned} \text{attn}(\mathbf{q}_t, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) &= \sum_i \frac{\exp(\mathbf{q}_t \cdot \mathbf{k}_i / \sigma^2)}{\sum_j \exp(\mathbf{q}_t \cdot \mathbf{k}_j / \sigma^2)} \mathbf{v}_i^\top \\ &\stackrel{(1)}{\approx} \sum_i \frac{\phi(\mathbf{q}_t)^\top \phi(\mathbf{k}_i) \mathbf{v}_i^\top}{\sum_j \phi(\mathbf{q}_t) \cdot \phi(\mathbf{k}_j)} \\ &= \frac{\phi(\mathbf{q}_t)^\top \sum_i \phi(\mathbf{k}_i) \otimes \mathbf{v}_i}{\phi(\mathbf{q}_t) \cdot \sum_j \phi(\mathbf{k}_j)} = \text{RFA}(\mathbf{q}_t, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) \end{aligned}$$



# Random feature attention [ICLR subm.]

- Exponential becomes Gaussian if  $x, y$  are normalized (Rawat et al., 2019)

$$\exp(x \cdot y / \sigma^2) = \exp(1/\sigma^2) \exp(-\|x - y\|^2 / 2\sigma^2)$$

temperature

- Therefore RFA is derived as follows

$$\begin{aligned} \text{attn}(\mathbf{q}_t, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) &= \sum_i \frac{\exp(\mathbf{q}_t \cdot \mathbf{k}_i / \sigma^2)}{\sum_j \exp(\mathbf{q}_t \cdot \mathbf{k}_j / \sigma^2)} \mathbf{v}_i^\top \\ &\stackrel{(1)}{\approx} \sum_i \frac{\phi(\mathbf{q}_t)^\top \phi(\mathbf{k}_i) \mathbf{v}_i^\top}{\sum_j \phi(\mathbf{q}_t) \cdot \phi(\mathbf{k}_j)} \\ &= \frac{\phi(\mathbf{q}_t)^\top \sum_i \phi(\mathbf{k}_i) \otimes \mathbf{v}_i}{\phi(\mathbf{q}_t) \cdot \sum_j \phi(\mathbf{k}_j)} = \text{RFA}(\mathbf{q}_t, \{\mathbf{k}_i\}, \{\mathbf{v}_i\}) \end{aligned}$$

Reparameterization trick

$$\tilde{\mathbf{w}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d),$$

$$\mathbf{w}_i = \sigma \circ \tilde{\mathbf{w}}_i.$$

(Kingma & Welling, 2014)

# RFA variants [ICLR subm.]

- **Non-causal:** we compute  $S, z$  only once for the whole sequence

$$\frac{\phi(\mathbf{q}_t)^\top \overbrace{\sum_i \phi(\mathbf{k}_i) \otimes \mathbf{v}_i}^S}{\phi(\mathbf{q}_t) \cdot \underbrace{\sum_j \phi(\mathbf{k}_j)}_z}$$

Encoder                      Cross

$\mathcal{O}(M)$                        $\mathcal{O}(M + N)$

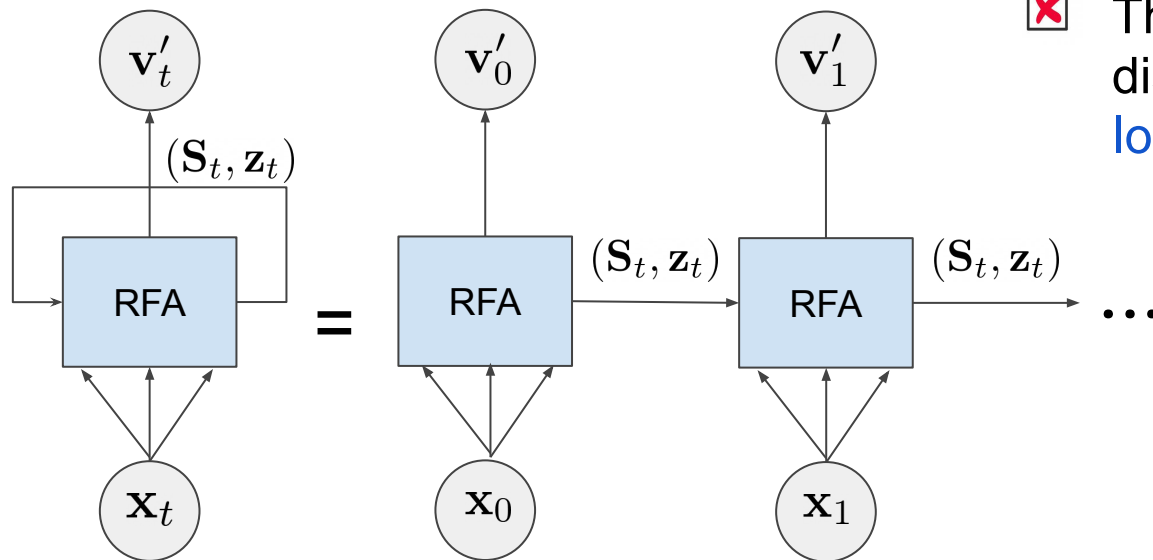
- **Causal:** we compute  $S_t$  and  $z_t$  iteratively at each step

$$\frac{\phi(\mathbf{q}_t)^\top \overbrace{\sum_{i \leq t} \phi(\mathbf{k}_i) \otimes \mathbf{v}_i}^{S_t}}{\phi(\mathbf{q}_t) \cdot \underbrace{\sum_{j \leq t} \phi(\mathbf{k}_j)}_{z_t}}$$

Decoder

$\mathcal{O}(N)$

# Recurrent formulation [ICLR subm.]

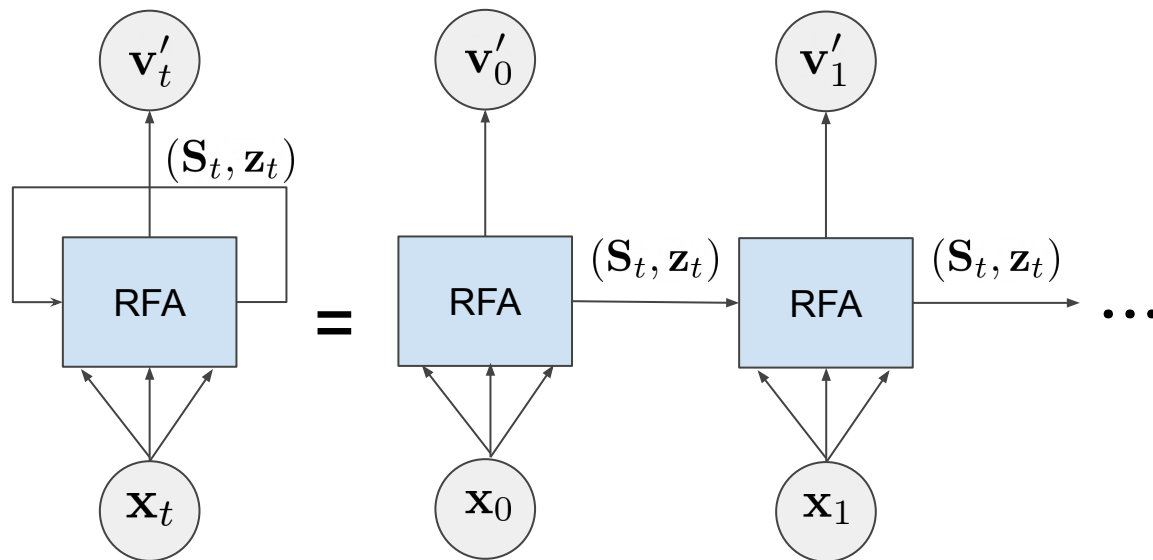


✗ There is no explicit modeling of distance or locality (Katharopoulos et al., 2020)

$$S_t = S_{t-1} + \phi(k_t) \otimes v_t$$

$$z_t = z_{t-1} + \phi(k_t)$$

# Gated-RFA: Learning with recency bias [ICLR subm.]



$$\begin{aligned} g_t &= \text{sigmoid}(\mathbf{w}_g \cdot \mathbf{x}_t + b_g), \\ \mathbf{S}_t &= g_t \mathbf{S}_{t-1} + (1 - g_t) \phi(\mathbf{k}_t) \otimes \mathbf{v}_t, \\ \mathbf{z}_t &= g_t \mathbf{z}_{t-1} + (1 - g_t) \phi(\mathbf{k}_t). \end{aligned}$$

# Machine translation

Model	WMT14		IWSLT14	Speed
	EN-DE	EN-FR	DE-EN	
BASE	28.1	39.0	34.6	1.0×
$\phi_{\text{elu}}$ (Katharopoulos et al., 2020)	21.3	34.0	29.9	2.0×
RFA-Gaussian	28.0	39.2	34.5	1.8×
RFA-arccos	28.1	38.9	34.4	1.9×
RFA-GATE-Gaussian	28.1	39.0	34.6	1.8×
RFA-GATE-arccos	28.2	39.2	34.4	1.9×

Higher is better

BLEU scores on MT.

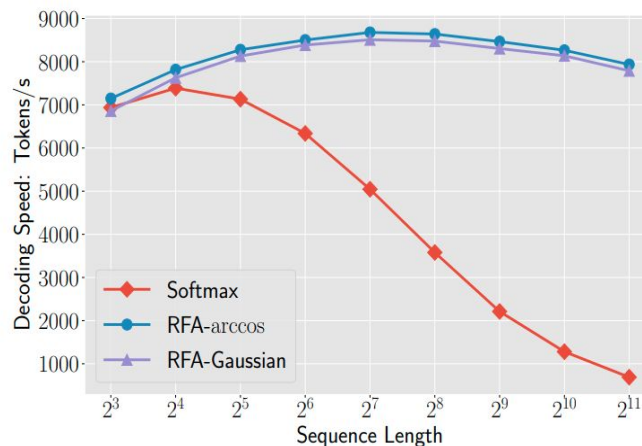
- Double the speed for short sequences with similar quality (2X speedup)
- Superiority to linear transformer shows the importance of feature map

# Language modeling

Model	Small		Big	
	Dev.	Test	Dev.	Test
BASE	33.0	34.5	24.5	26.2
$\phi_{\text{elu}}$ (Katharopoulos et al., 2020)	38.4	40.1	28.7	30.2
RFA-Gaussian	33.6	35.7	25.8	27.5
RFA-arccos	36.0	37.7	26.4	28.1
RFA-GATE-Gaussian	<b>31.3</b>	<b>32.7</b>	<b>23.2</b>	<b>25.0</b>
RFA-GATE-arccos	<b>32.8</b>	<b>34.0</b>	24.8	26.3
RFA-GATE-Gaussian-Stateful	<b>29.4</b>	<b>30.5</b>	<b>22.0</b>	<b>23.5</b>

Lower is better

Perplexity on Wikitext-103.

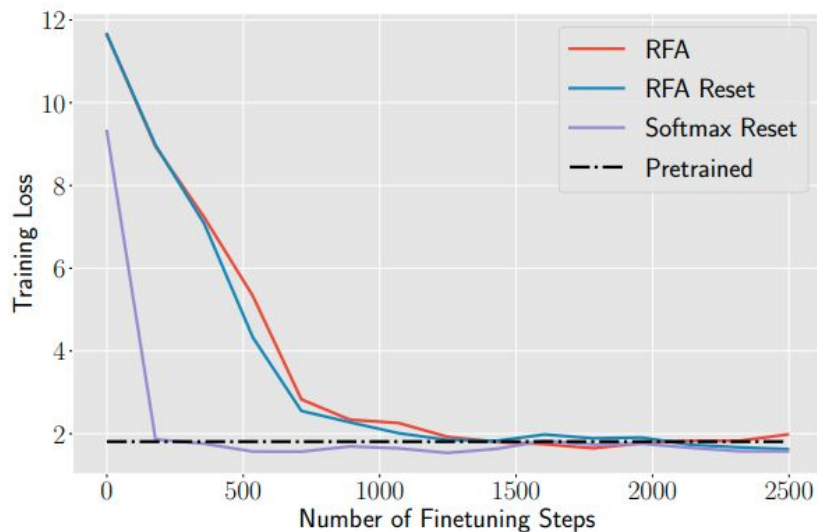


Higher is better

Decoding speed

- RFA-gate is better than baseline with up to 17X decoding speed
- Competitive speed-quality tradeoff in the long range arena (5X speedup)

# “Streamlining” pretrained models

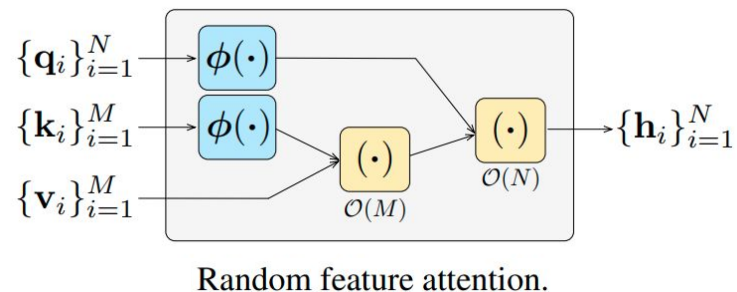


Finetuning RFA from a pretrained softmax model.

- Pretrained softmax parameters are as useful as random ones
- RFA can recover the pretraining loss with a few iterations
- Potential to reduce finetuning cost for large models (GPT3)

# Takeaways [ICLR subm.]

- General component with linear complexity for attending sequences
- Competitive trade-offs on both long and moderate length sequences
- New scalable attention variant that learns with recency bias





# Acknowledgments



L. Miculicich



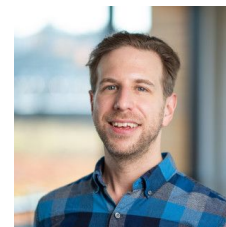
D. Ram



J. Henderson



A. Popescu-Belis



N. A. Smith



H. Peng



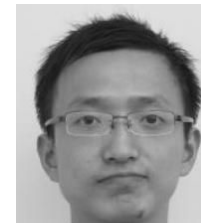
P. Mulcaire



D. Yogatama



R. Schwartz



L. Kong



Horizon 2020  
European Union Funding  
for Research & Innovation

Thank you!