

# Learning How to Say It: Language Generation post Deep Learning

Alexander M Rush

# Machine Learning for Multiclass Classification

$x$



# Machine Learning for Multiclass Classification

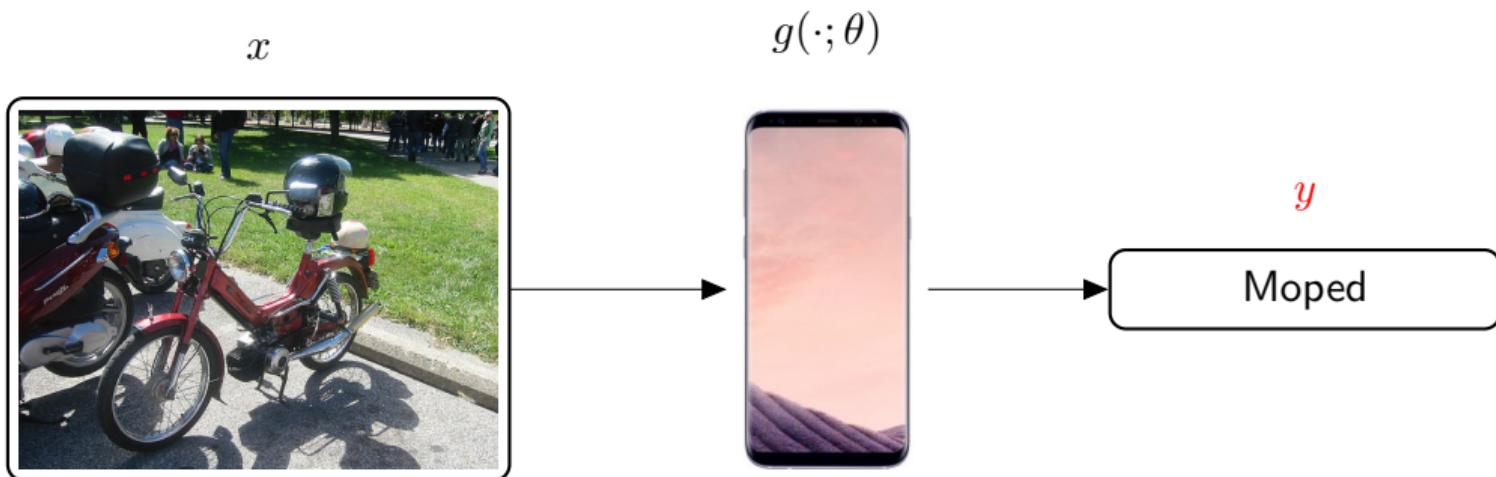
$x$



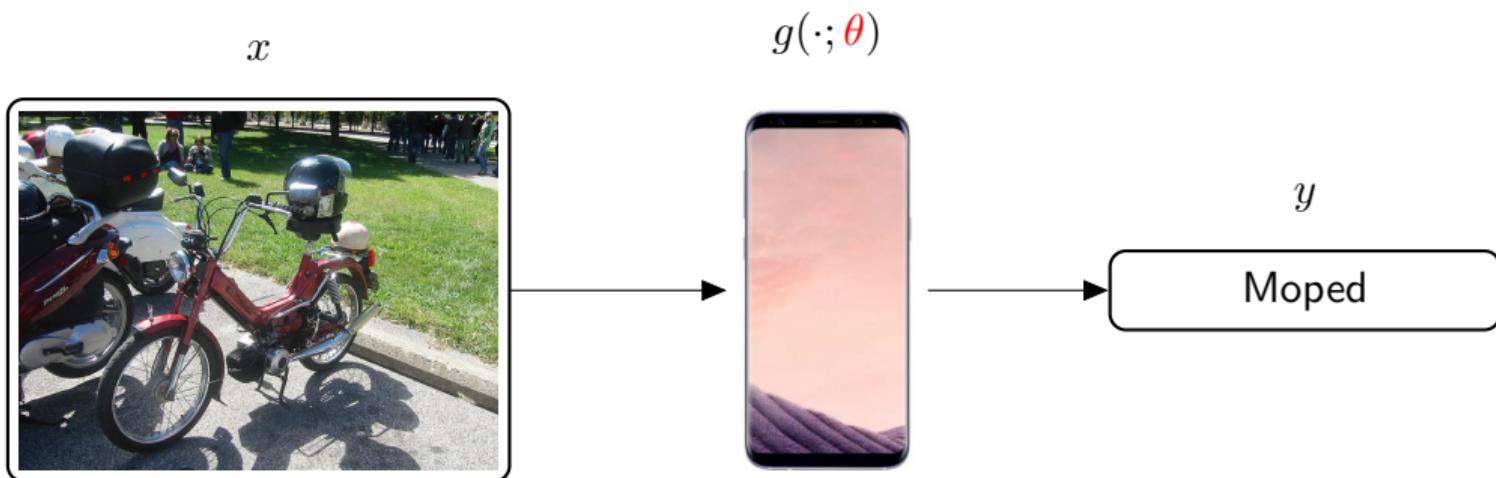
$g(\cdot; \theta)$



# Machine Learning for Multiclass Classification



# Machine Learning for Multiclass Classification



# Machine Learning for Text Generation

$$y_{1:T}^* = \arg \max_{y_{1:T}} f(y_{1:T}, x; \theta)$$

# Machine Learning for Text Generation

$$y_{1:T}^* = \arg \max_{y_{1:T}} f(y_{1:T}, \textcolor{red}{x}; \theta)$$

- Input  $\textcolor{red}{x}$ , what to talk about

# Machine Learning for Text Generation

$$y_{1:T}^* = \arg \max_{y_{1:T}} f(y_{1:T}, x; \theta)$$

- Input  $x$ , *what to talk about*
- Output text  $y_{1:T}^*$ , *how to say it*

# Machine Learning for Text Generation

$$y_{1:T}^* = \arg \max_{y_{1:T}} \textcolor{red}{f}(y_{1:T}, x; \theta)$$

- Input  $x$ , *what to talk about*
- Output text  $y_{1:T}^*$ , *how to say it*
- Model  $\textcolor{red}{f}(\cdot; \theta)$ , learned from data

# Part 1: Generating Text

Applications

# Machine Translation

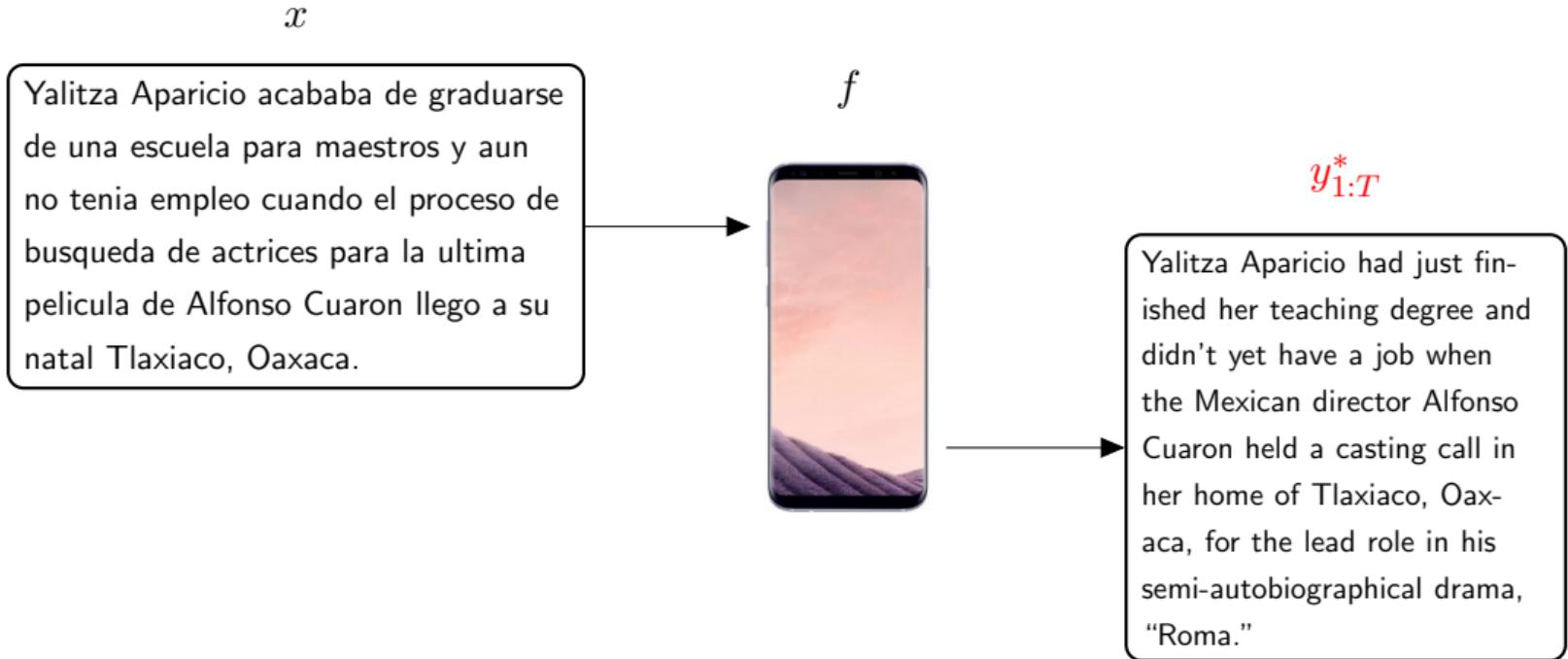
*x*

Yalitza Aparicio acababa de graduarse  
de una escuela para maestros y aun  
no tenia empleo cuando el proceso de  
busqueda de actrices para la ultima  
pelicula de Alfonso Cuaron llego a su  
natal Tlaxiaco, Oaxaca.

*f*



# Machine Translation



## Translation Performance

### Evaluation Metric:

Target: [Yalitza Aparicio had] just [finished her] teaching [degree] .

Predict: [Yalitza Aparicio had] recently [finished her] [degree] .

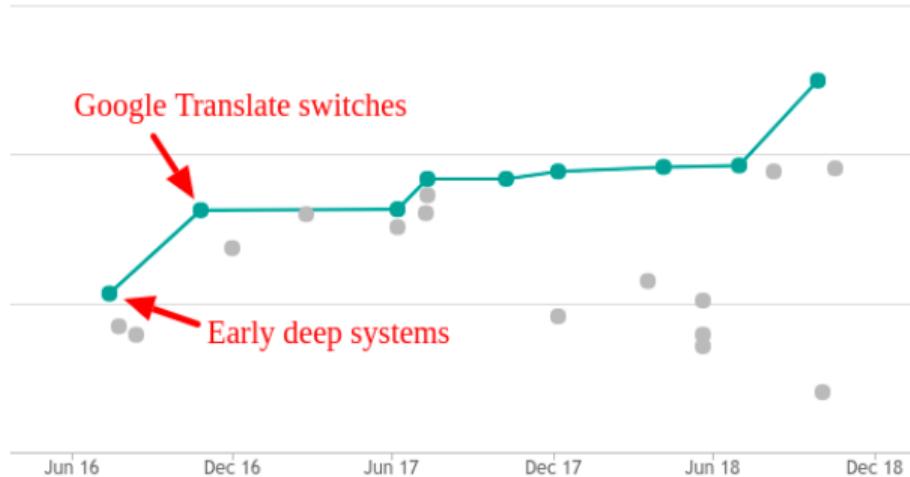
# Translation Performance

## Evaluation Metric:

Target: [Yalitza Aparicio had] just [finished her] teaching [degree] .

Predict: [Yalitza Aparicio had] recently [finished her] [degree] .

## Deep Learning Performance:



# Sentence Summarization

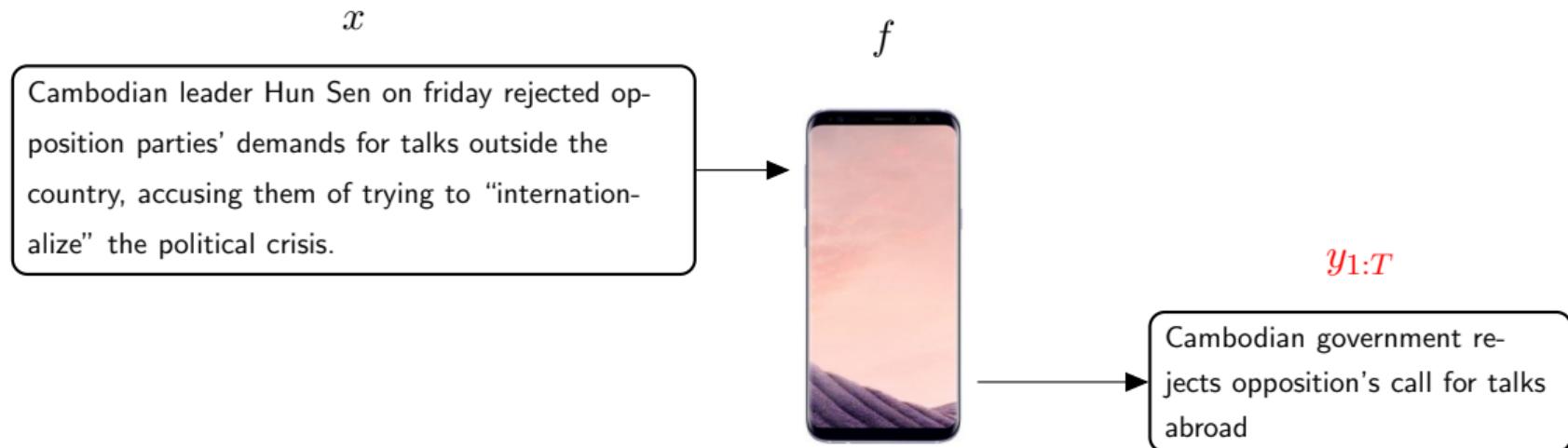
$x$

Cambodian leader Hun Sen on Friday rejected opposition parties' demands for talks outside the country, accusing them of trying to "internationalize" the political crisis.

$f$



# Sentence Summarization



# GigaWord Dataset

(Rush et al. [2015] w/ Facebook)

Sep 13, 3:17 PM EDT

## GERMANY IMPLEMENTS TEMPORARY BORDER CHECKS TO LIMIT MIGRANTS

BY GEIR MOULSON AND SHAWN POGATCHNIK  
ASSOCIATED PRESS

0

BERLIN (AP) -- Germany introduced temporary border controls Sunday to stem the tide of thousands of refugees streaming across its frontier, sending a clear message to its European partners that it needs more help with an influx that is straining its ability to cope.

Germany is a preferred destination for many people fleeing Syria's civil war and other troubled nations in the migration crisis that has bitterly divided Europe. They have braved dangerous sea crossings in flimsy



AP Photo/Kay Nietfeld

- Several million headlines paired with article leads.
- Simple model for abstractive summarization / compression.
- Benchmark dataset for early deep summarization work

# Sentence Summarization



# Document Summary

London, England (reuters) – Harry Potter star Daniel Radcliffe gains access to a reported \$20 million fortune as he turns 18 on monday, but he insists the money won't cast a spell on him. Daniel Radcliffe as harry potter in "Harry Potter and the Order of the Phoenix" to the disappointment of gossip columnists around the world , the young actor says he has no plans to fritter his cash away on fast cars , drink and celebrity parties . " i do n't plan to be one of those people who , as soon as they turn 18 , suddenly buy themselves a massive sports car collection or something similar , " he told an australian interviewer earlier this month . " i do n't think i 'll be particularly extravagant " . " the things i like buying are things that cost about 10 pounds – books and cds and dvds . " at 18 , radcliffe will be able to gamble in a casino , buy a drink in a pub or see the horror film " hostel : part ii , " currently six places below his number one movie on the uk box office chart . details of how he 'll mark his landmark birthday are under wraps . his agent and publicist had no comment on his plans . " i 'll definitely have some sort of party , " he said in an interview ...



# Document Summary

London, England (reuters) – Harry Potter star Daniel Radcliffe gains access to a reported \$20 million fortune as he turns 18 on monday, but he insists the money won't cast a spell on him. Daniel Radcliffe as harry potter in "Harry Potter and the Order of the Phoenix" to the disappointment of gossip columnists around the world , the young actor says he has no plans to fritter his cash away on fast cars , drink and celebrity parties . " i do n't plan to be one of those people who , as soon as they turn 18 , suddenly buy themselves a massive sports car collection or something similar , " he told an australian interviewer earlier this month . " i do n't think i 'll be particularly extravagant " . " the things i like buying are things that cost about 10 pounds – books and cds and dvds . " at 18 , radcliffe will be able to gamble in a casino , buy a drink in a pub or see the horror film " hostel : part ii , " currently six places below his number one movie on the uk box office chart . details of how he 'll mark his landmark birthday are under wraps . his agent and publicist had no comment on his plans . " i 'll definitely have some sort of party , " he said in an interview ...



Harry Potter star Daniel Radcliffe gets \$20m fortune as he turns 18 monday. Young actor says he has no plans to fritter his cash away. Radcliffe 's earnings from first five potter films have been held in trust fund.

# Document Summarization



# Talk about Data

(Wiseman et al. [2017a])

TEAM	WIN	LOSS	PTS	FG_PCT	RB	AS ...
Heat	11	12	103	49	47	27
Hawks	7	15	95	43	33	20

PLAYER	AS	RB	PT	FG	FGA	CITY ...
Tyler Johnson	5	2	27	8	16	Miami
Dwight Howard	11	17	23	9	11	Atlanta
Paul Millsap	2	9	21	8	12	Atlanta
Goran Dragic	4	2	21	8	17	Miami
Wayne Ellington	2	3	19	7	15	Miami
Dennis Schroder	7	4	17	8	15	Atlanta
Rodney McGruder	5	5	11	3	8	Miami
...						



# Talk about Data

(Wiseman et al. [2017a])

TEAM	WIN	LOSS	PTS	FG_PCT	RB	AS ...
Heat	11	12	103	49	47	27
Hawks	7	15	95	43	33	20
PLAYER	AS	RB	PT	FG	FGA	CITY ...
Tyler Johnson	5	2	27	8	16	Miami
Dwight Howard	11	17	23	9	11	Atlanta
Paul Millsap	2	9	21	8	12	Atlanta
Goran Dragic	4	2	21	8	17	Miami
Wayne Ellington	2	3	19	7	15	Miami
Dennis Schroder	7	4	17	8	15	Atlanta
Rodney McGruder	5	5	11	3	8	Miami
...						



The Atlanta Hawks defeated the Miami Heat, 103 - 95, at Philips Arena on Wednesday. Atlanta was in desperate need of a win and they were able to take care of a short-handed Miami team here. Defense was key for the Hawks, as they held the Heat to 42 percent shooting and forced them to commit 16 turnovers. Atlanta also dominated in the paint, winning the rebounding battle, 47 - 34, and outscoring them in the paint 58 - 26. The Hawks shot 49 percent from the field and assisted on 27 of their 43 made baskets. This was a near wire-to-wire win for the Hawks, as Miami held just one lead in the first five minutes. Miami ( 7 - 15 ) are as beat-up as anyone right now and it's taking a toll on the heavily used starters. Hassan Whiteside really struggled in this game, as he amassed eight points, 12 rebounds and one blocks on 4 - of - 12 shooting ...

# E2E Challenge 2018

---

<b>MR</b>	name[The Golden Palace], eatType[coffee shop], food[Fast food], priceRange[cheap], customer rating[5 out of 5], area[riverside]
<b>Reference</b>	A coffee shop located on the riverside called The Golden Palace, has a 5 out of 5 customer rating. Its price range are fairly cheap for its excellent Fast food.

---

Submitter	Affiliation	System name	P?	BLEU	NIST	METEOR	ROUGE_L▲	CIDEr
HarvardNLP & Henry Elder	Harvard SEAS & Adapt	main_1_support_3		0.6737	8.6061	0.4523	0.7084	2.3056
Biao Zhang	Xiamen University	bzhang_submit	✓	0.6545	8.1840	0.4392	0.7083	2.1012
HarvardNLP & Henry Elder	Harvard SEAS & Adapt	main_1_support_2		0.6618	8.6025	0.4571	0.7038	2.3371
Shubham Agarwal	NLE	submission_third		0.6676	8.5416	0.4485	0.6991	2.2276
Shubham Agarwal	NLE	submission_second		0.6669	8.5388	0.4484	0.6991	2.2239
Thomson Reuters NLP	Thomson Reuters	NonPrimary_4_test_output_beam_5_model_13_post		0.6742	8.6590	0.4499	0.6983	2.3018

# Talk about the Diagrams (Deng et al. [2016] w/ Bloomberg)

$$\mathcal{K}^L(\sigma = 2) = \begin{pmatrix} -\frac{d^2}{dx^2} + 4 - \frac{3}{\cosh^2 x} & \frac{3}{\cosh^2 x} \\ \frac{3}{\cosh^2 x} & -\frac{d^2}{dx^2} + 4 - \frac{3}{\cosh^2 x} \end{pmatrix},$$



# Talk about the Diagrams (Deng et al. [2016] w/ Bloomberg)

$$\mathcal{K}^L(\sigma = 2) = \begin{pmatrix} -\frac{d^2}{dx^2} + 4 - \frac{3}{\cosh^2 x} & \frac{3}{\cosh^2 x} \\ \frac{3}{\cosh^2 x} & -\frac{d^2}{dx^2} + 4 - \frac{3}{\cosh^2 x} \end{pmatrix},$$



```
{ \cal K } ^ { L } ( \sigma = 2 ) = \left( \begin{array}{cc} - \frac{d^2}{dx^2} + 4 - \frac{3}{\operatorname{cosh}^2 x} & \frac{3}{\operatorname{cosh}^2 x} \\ \frac{3}{\operatorname{cosh}^2 x} & - \frac{d^2}{dx^2} + 4 - \frac{3}{\operatorname{cosh}^2 x} \end{array} \right) \quad ,
```

$$A_0^3(\alpha' \rightarrow 0) = 2g_d\, \varepsilon_\lambda^{(1)} \varepsilon_\mu^{(2)} \varepsilon_\nu^{(3)} \left\{ \eta^{\lambda\mu} (p_1^\nu - p_2^\nu) + \eta^{\lambda\nu} (p_3^\mu - p_1^\mu) + \eta^{\mu\nu} (p_2^\lambda - p_3^\lambda) \right\}.$$

$$(\Lambda_{-0})^3(\alpha'\prime)\rightarrow 0)=2g_{-d},\backslash\varepsilon_\lambda^{(1)}\varepsilon_\mu^{(2)}\varepsilon_\nu^{(3)}\left\{\eta^{\lambda\mu}(p_1^\nu-p_2^\nu)+\eta^{\lambda\nu}(p_3^\mu-p_1^\mu)+\eta^{\mu\nu}(p_2^\lambda-p_3^\lambda)\right\}.$$

$$(\Lambda_{-\mu})^3(\nu)\left(\eta^{\lambda\mu}(p_{-1}^\nu-p_{-2}^\nu)+\eta^{\lambda\nu}(p_{-3}^\mu-p_{-1}^\mu)+\eta^{\mu\nu}(p_{-2}^\lambda-p_{-3}^\lambda)\right).$$

\label{17}

$$\begin{cases} \delta_\epsilon B & \sim \epsilon F, \\ \delta_\epsilon F & \sim \partial \epsilon + \epsilon B, \end{cases}$$

$$\left.\left(\begin{array}{ccl}\delta_\epsilon B & \sim & \epsilon F,\\ \delta_\epsilon F & \sim & \partial \epsilon + \epsilon B,\end{array}\right.\right.$$

$$\int\limits_{\mathcal{L}_{d-1}^d}f(H)d\nu_{d-1}(H)=c_3\int\limits_{\mathcal{L}_2^A}\int\limits_{\mathcal{L}_{d-1}^L}f(H)[H,A]^2d\nu_{d-1}^L(H)d\nu_2^A(L).$$

$$\int \limits_{\{\mathcal{L}\}} \{\mathcal{L}\}^{(d-1)} f(H) d\nu_{(d-1)}(H) = c_{-3} \int \limits_{\{\mathcal{L}\}} \{\mathcal{L}\}^{(A)} \int \limits_{\{\mathcal{L}\}} \{\mathcal{L}\}^{(d-1)} f(H) [H,A]^{(2)} d\nu_{(d-1)}(L) d\nu_{(A)}(L).$$

$$J=\left(\begin{array}{cc}\alpha^t&\tilde{f}_2\\f_1&\tilde{A}\end{array}\right)\left(\begin{array}{cc}0&0\\0&L\end{array}\right)\left(\begin{array}{cc}\alpha&\tilde{f}_1\\f_2&A\end{array}\right)=\left(\begin{array}{cc}\tilde{f}_2Lf_2&\tilde{f}_2LA\\\tilde{A}Lf_2&\tilde{A}LA\end{array}\right)$$

$$J=\left(\begin{array}{cc}\alpha^t&\tilde{f}_2\\f_1&\tilde{A}\end{array}\right)\left(\begin{array}{cc}0&0\\0&L\end{array}\right)\left(\begin{array}{cc}\alpha&\tilde{f}_1\\f_2&A\end{array}\right)=\left(\begin{array}{cc}\tilde{f}_2Lf_2&\tilde{f}_2LA\\\tilde{A}Lf_2&\tilde{A}LA\end{array}\right)$$

$$\lambda_{n,1}^{(2)}=\frac{\partial \overline{H}_0}{\partial q_{n,0}}\; , lambda_{n,j_n}^{(2)}=\frac{\partial \overline{H}_0}{\partial q_{n,j_n-1}}-\mu_{n,j_n-1}\; , \;\; j_n=2,3,\cdots,m_n-1\; .$$

$$\lambda_{n,1}^{(2)}=\frac{\partial \overline{H}_0}{\partial q_{n,0}}\; , lambda_{n,j_n}^{(2)}=\frac{\partial \overline{H}_0}{\partial q_{n,j_n-1}}-\mu_{n,j_n-1}\; , \;\; j_n=2,3,\cdots,m_n-1\; .$$

$$(P_{ll'}-K_{ll'})\phi'(z_q)|\chi>=0$$

$$(P_{\{ll'\}}-K_{\{ll'\}})\phi'(z_{\{q\}})|\chi>=0$$



# 1 Introduction

# Machine Learning for Natural Language

What types of models get used for this equation?

$$\arg \max_{y_{1:T}} \textcolor{red}{f}(y_{1:T}, x; \theta)$$

# State-of-the-Art Natural Language Processing, circa 2009

**Task**  $(x, y)$

Syntax

Surface Structure

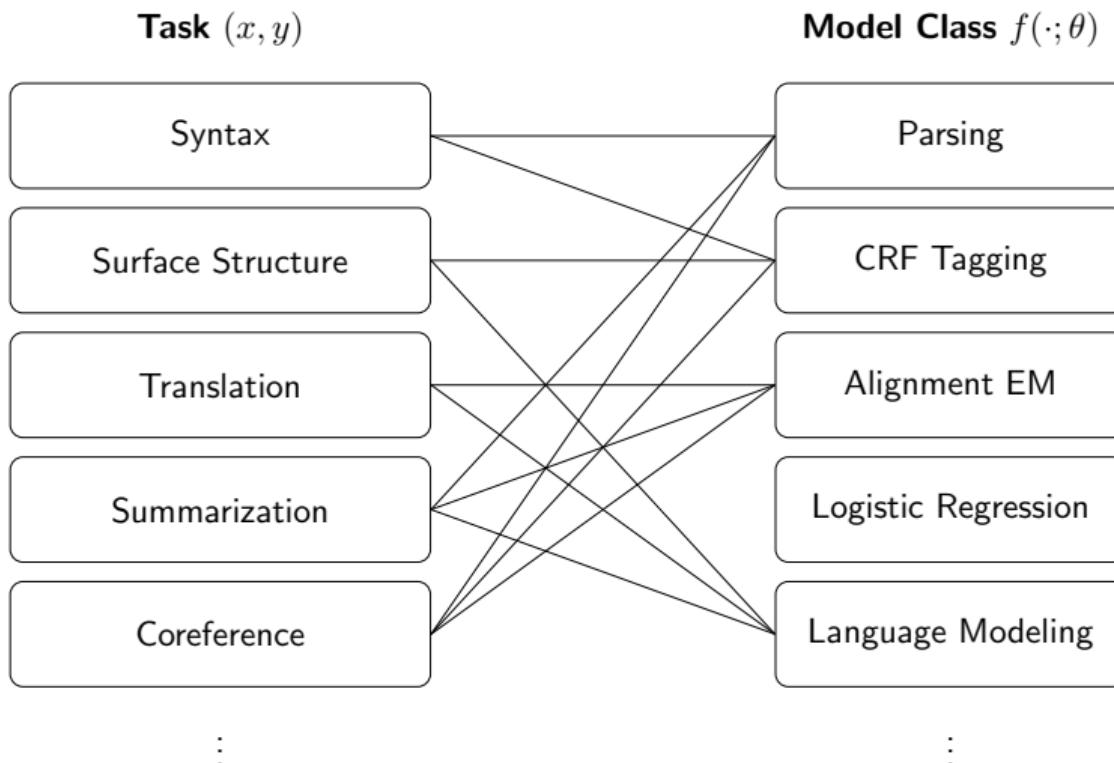
Translation

Summarization

Coreference

:

# State-of-the-Art Natural Language Processing, circa 2009



# State-of-the-Art Natural Language Processing, circa 2019

**Task** ( $x, y$ )

Syntax

Surface Structure

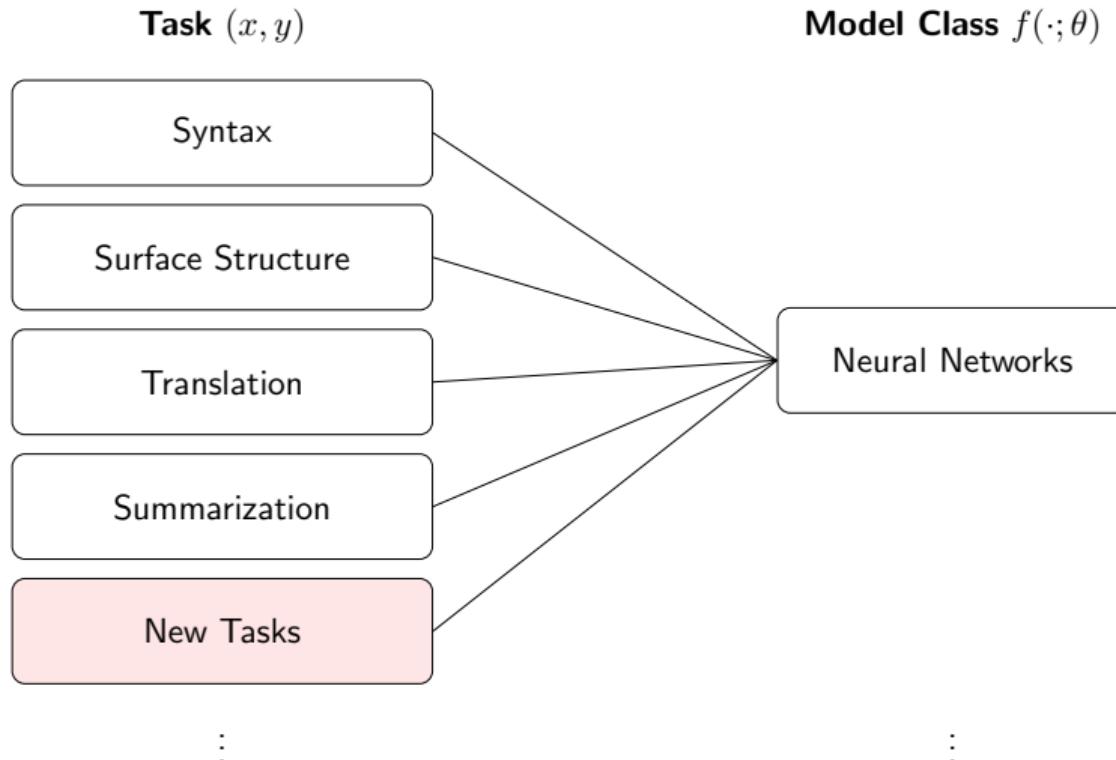
Translation

Summarization

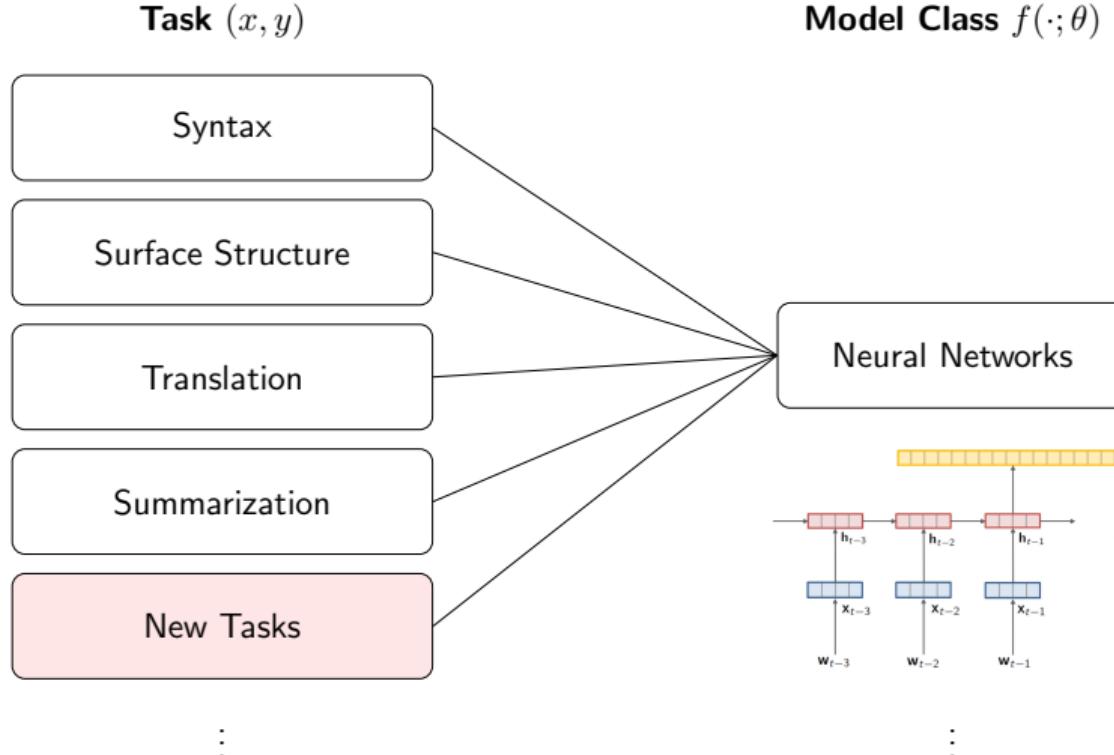
New Tasks

:

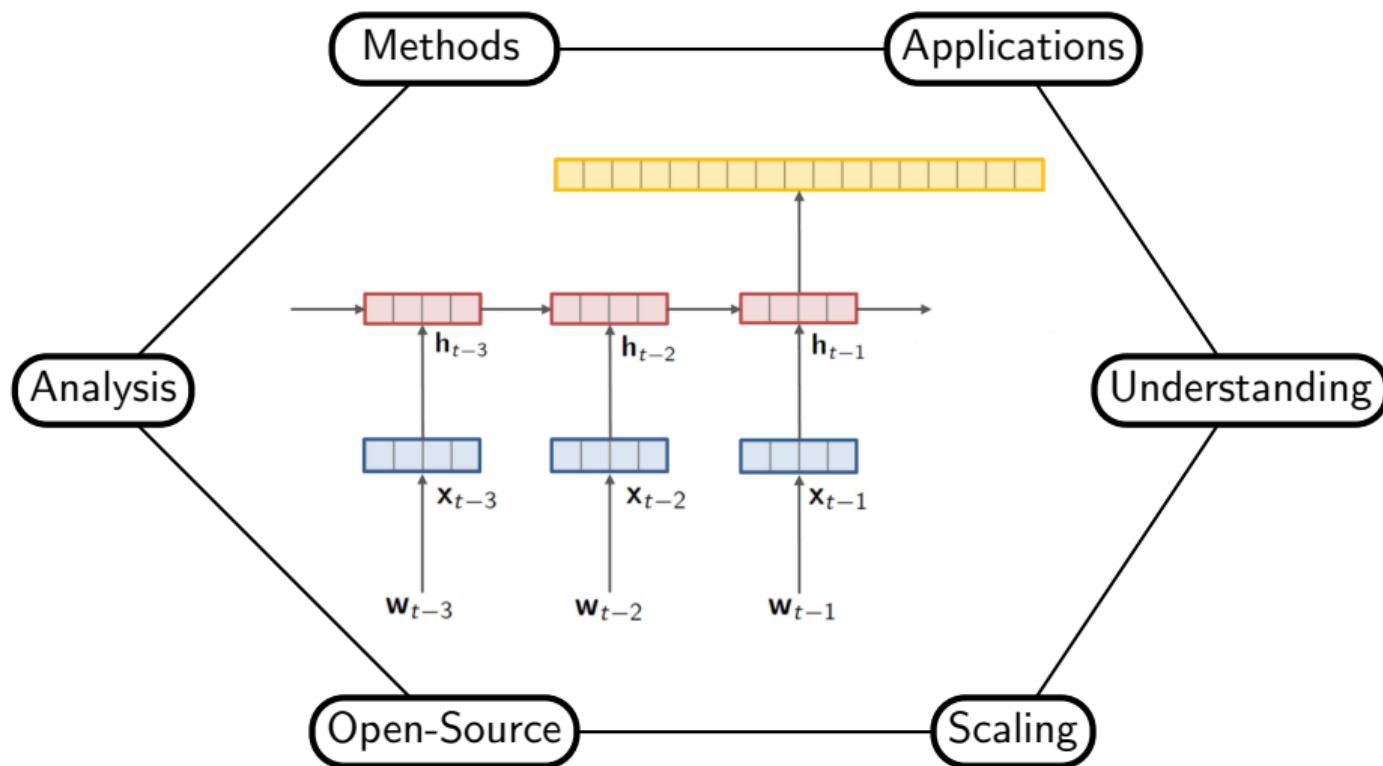
# State-of-the-Art Natural Language Processing, circa 2019



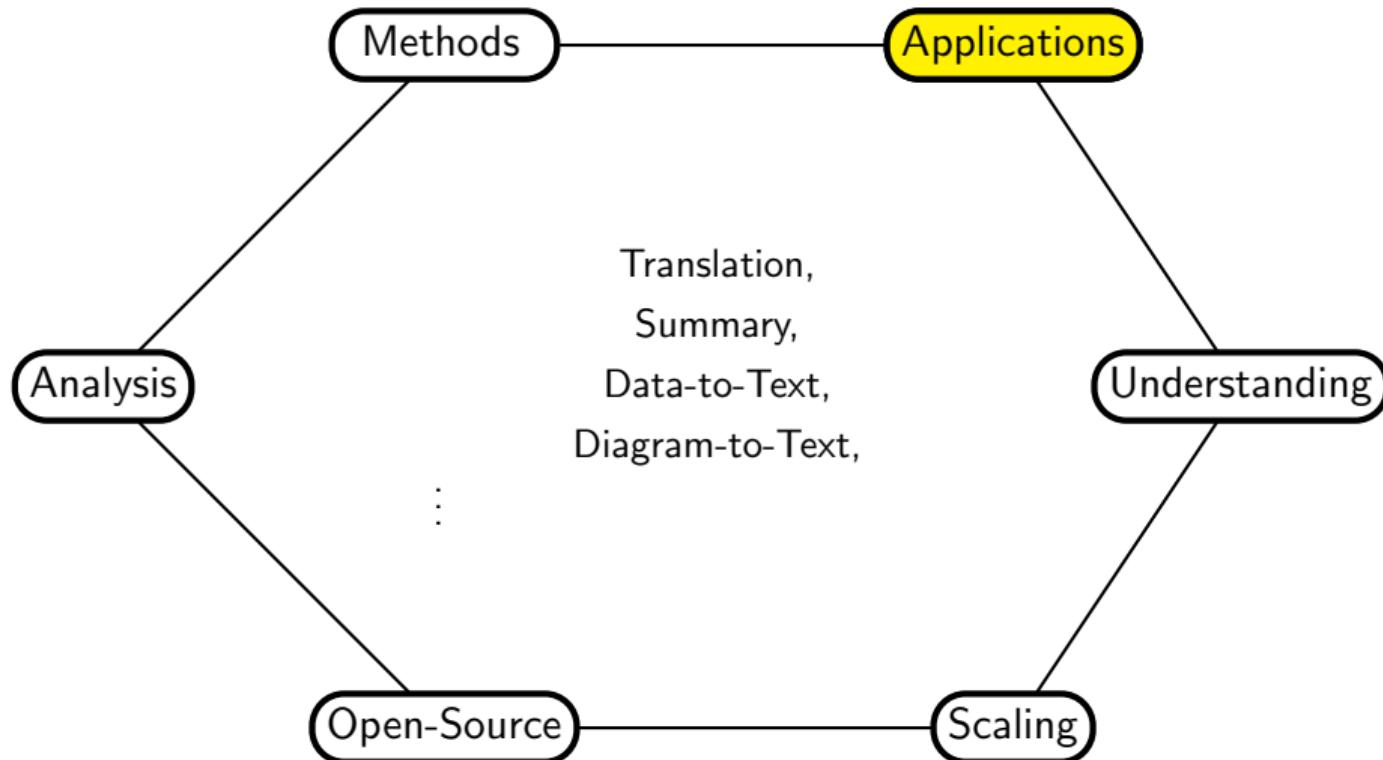
# State-of-the-Art Natural Language Processing, circa 2019



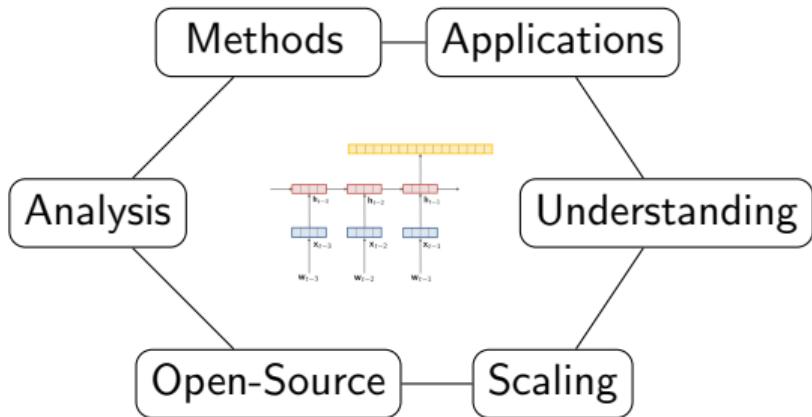
# Harvard NLP Deep Learning Research



# Harvard NLP Deep Learning Research

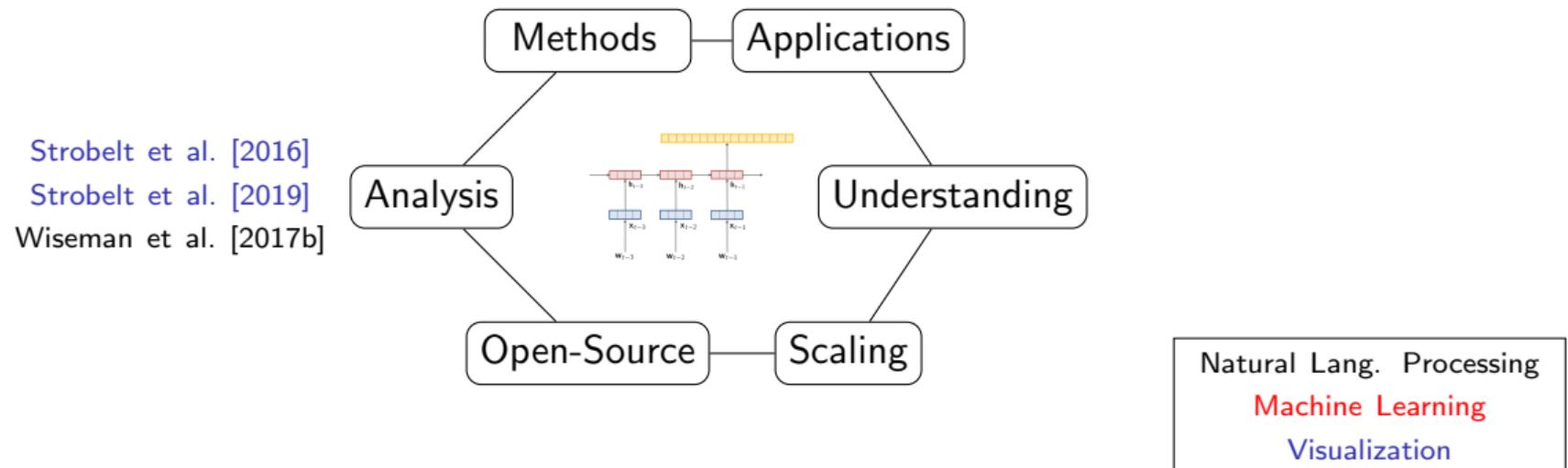


# Selected Harvard NLP Deep Learning Research



Natural Lang. Processing  
Machine Learning  
Visualization

# Selected Harvard NLP Deep Learning Research

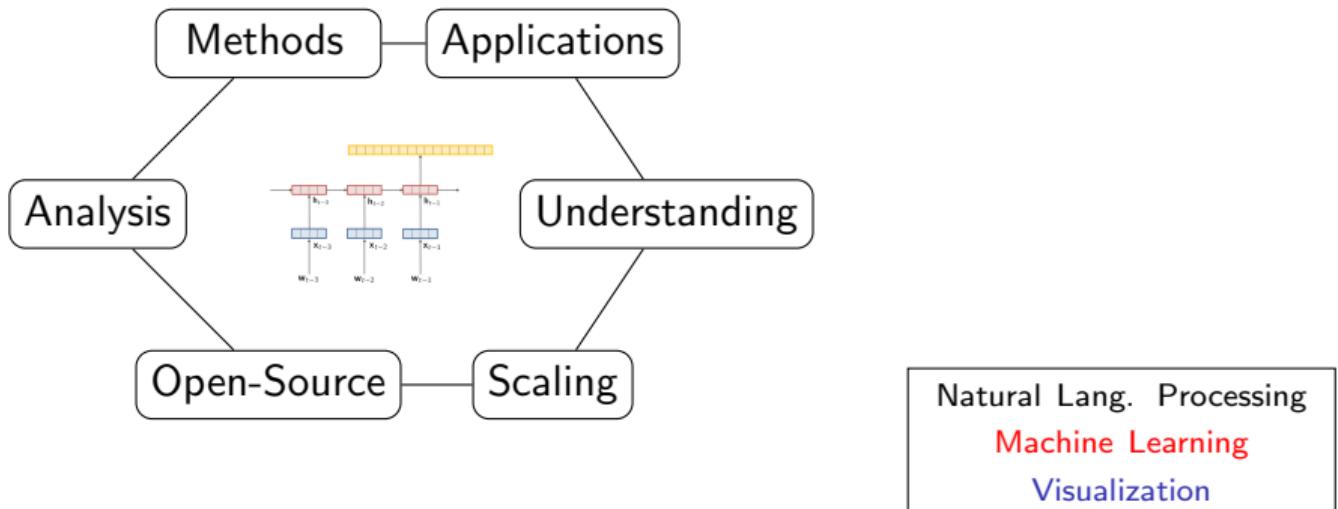


# Selected Harvard NLP Deep Learning Research

Kim et al. [2016]

Kim et al. [2017]

Wiseman et al. [2017a]

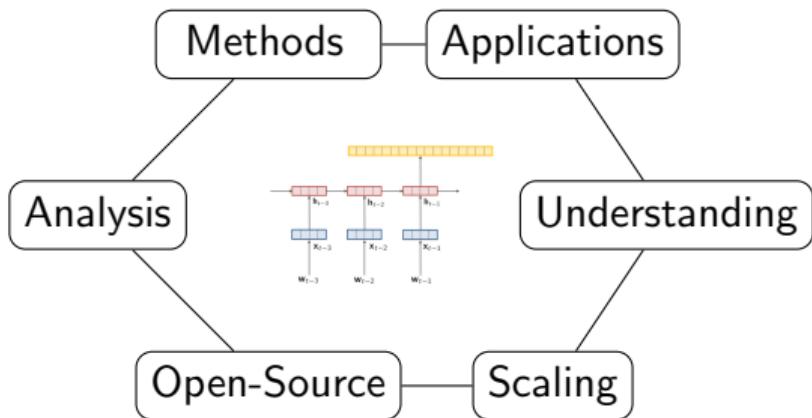


# Selected Harvard NLP Deep Learning Research

Rush et al. [2015]

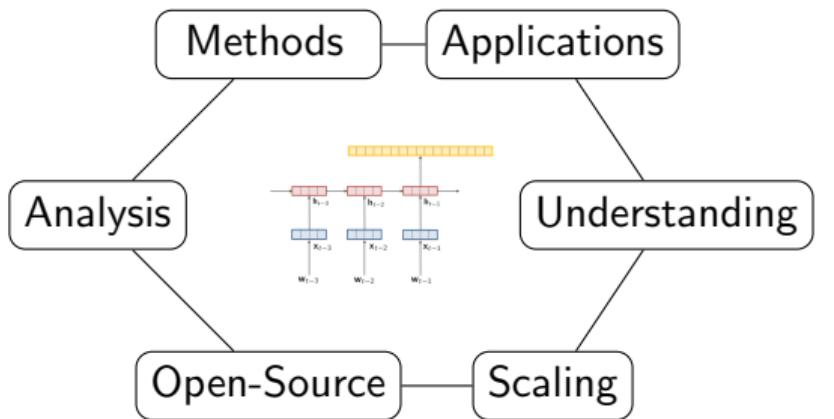
Deng et al. [2016]

Schmaltz et al. [2016]



Natural Lang. Processing  
Machine Learning  
Visualization

# Selected Harvard NLP Deep Learning Research



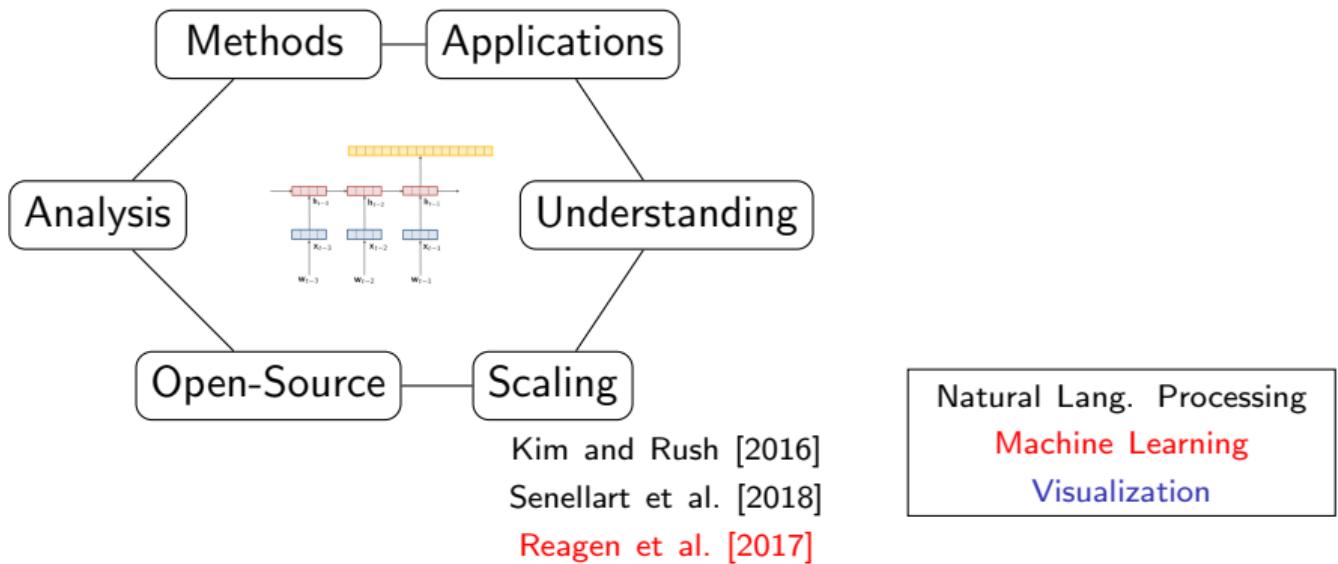
Wiseman et al. [2018]

Deng et al. [2018]

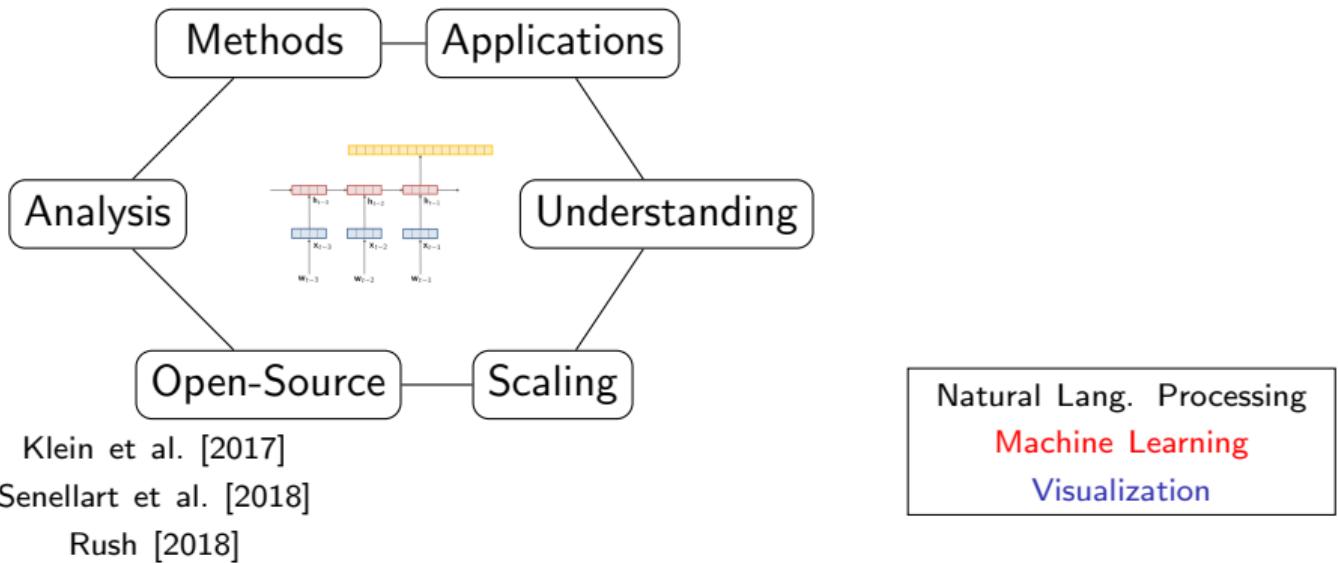
Kim et al. [2018]

Natural Lang. Processing  
Machine Learning  
Visualization

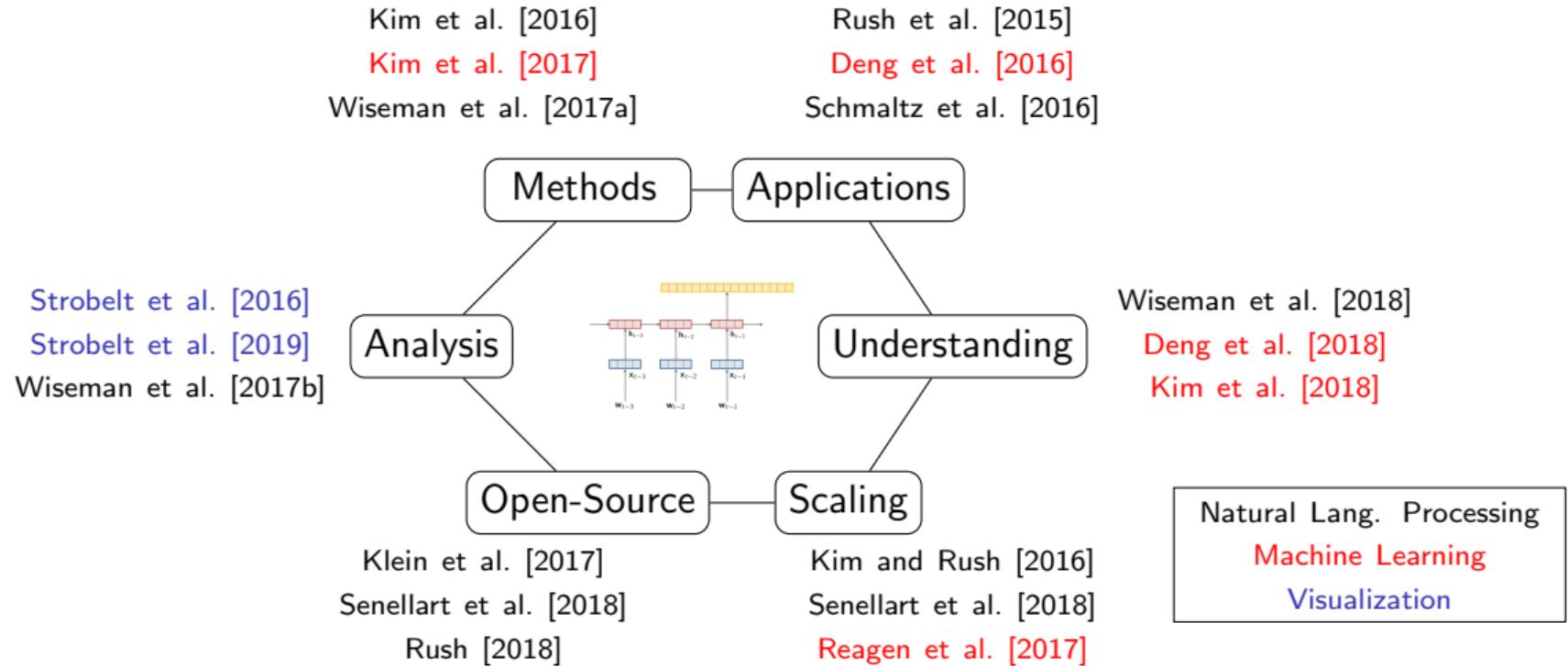
# Selected Harvard NLP Deep Learning Research



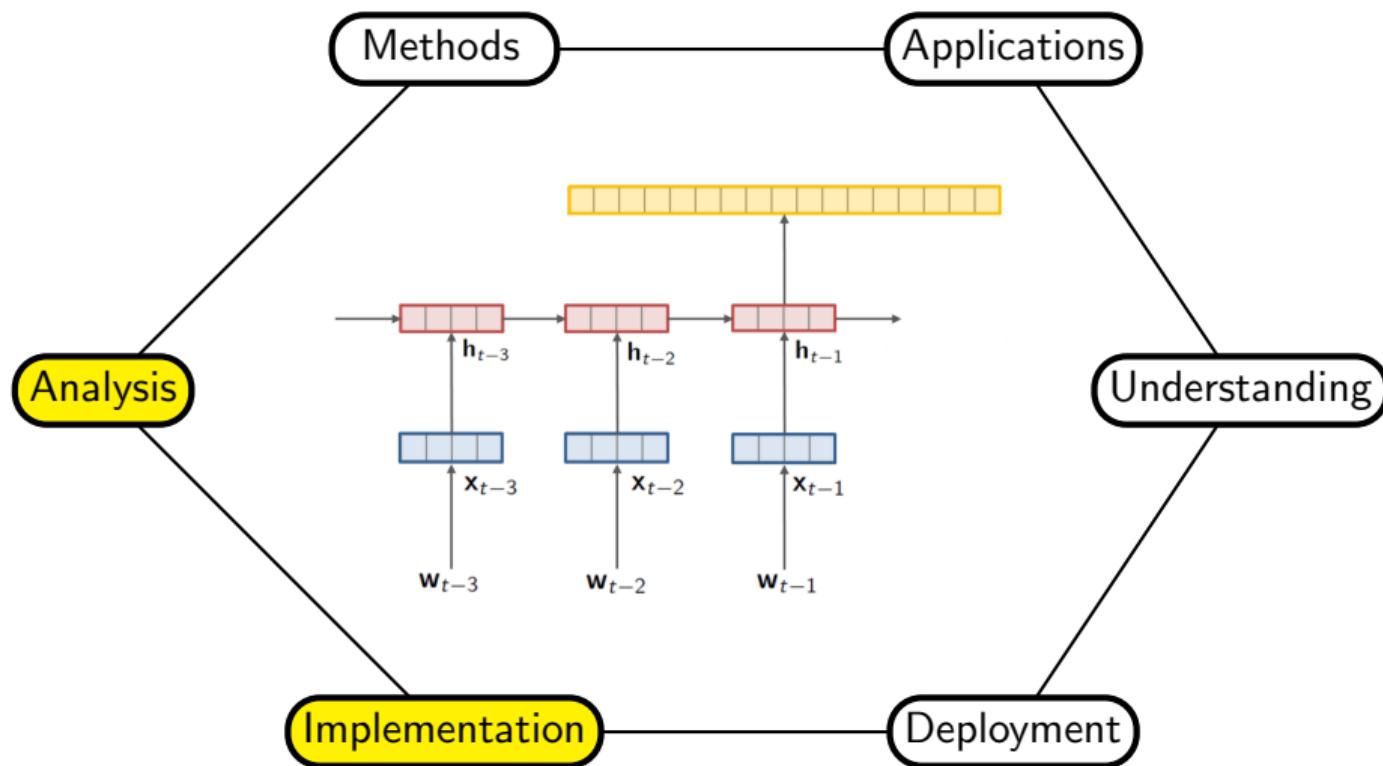
# Selected Harvard NLP Deep Learning Research



# Selected Harvard NLP Deep Learning Research



## Part 2: Deep Learning Internals



# Machine Learning for Text Generation

$$y_{1:T}^* = \arg \max_{y_{1:T}} \textcolor{red}{f}(y_{1:T}; x, \theta)$$

- Input  $x_{1:S}$ , *what to talk about*
- Output text  $y_{1:T}^*$ , *how to say it*
- Model  $\textcolor{red}{f}(\cdot; \theta)$ , learned from data

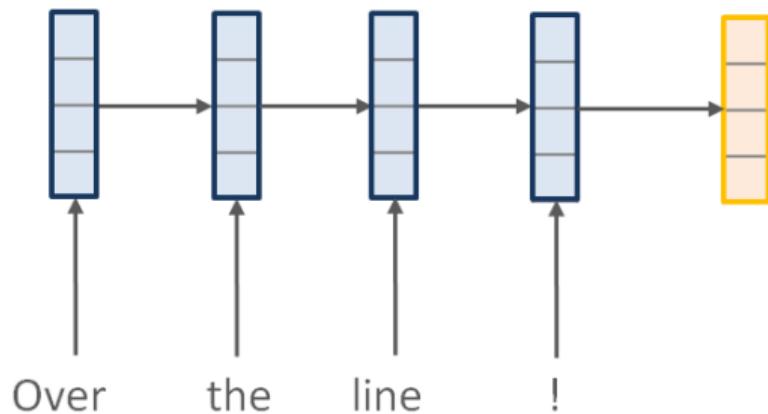
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$

Over the line !

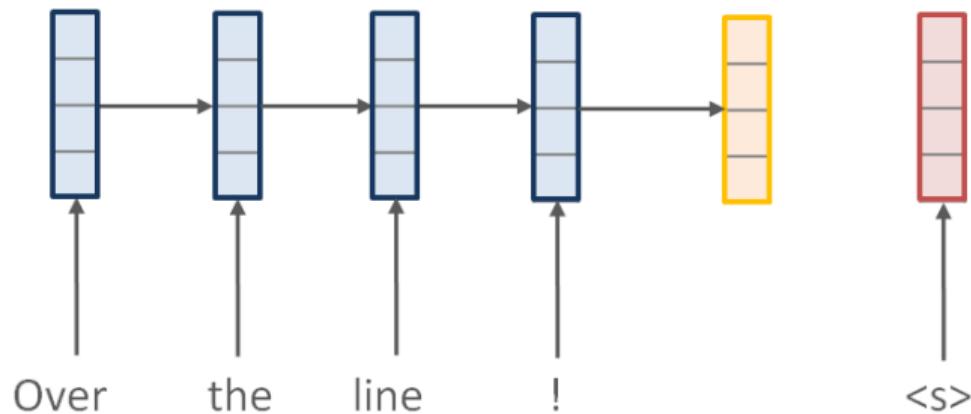
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



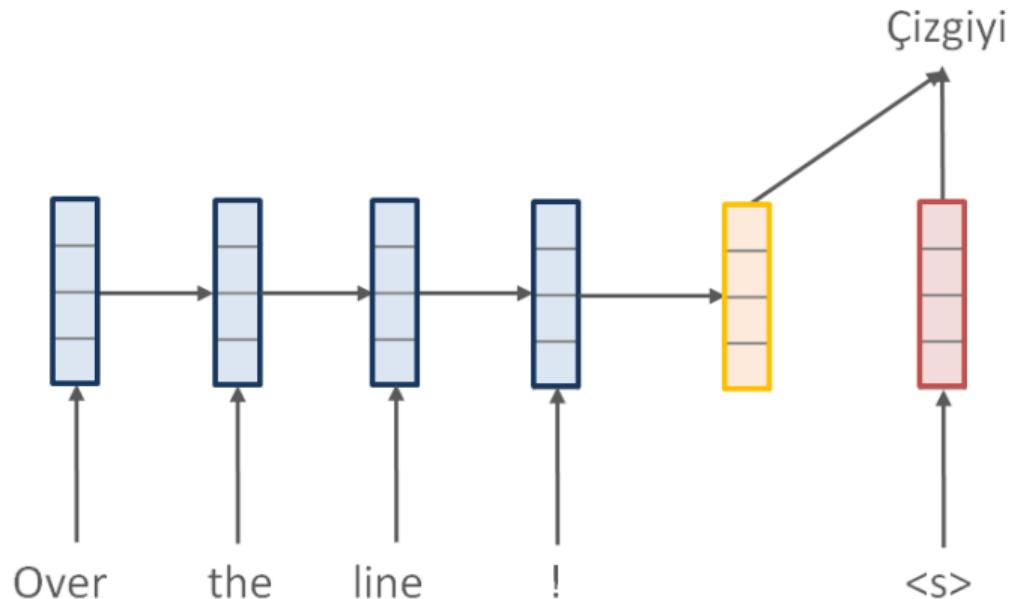
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



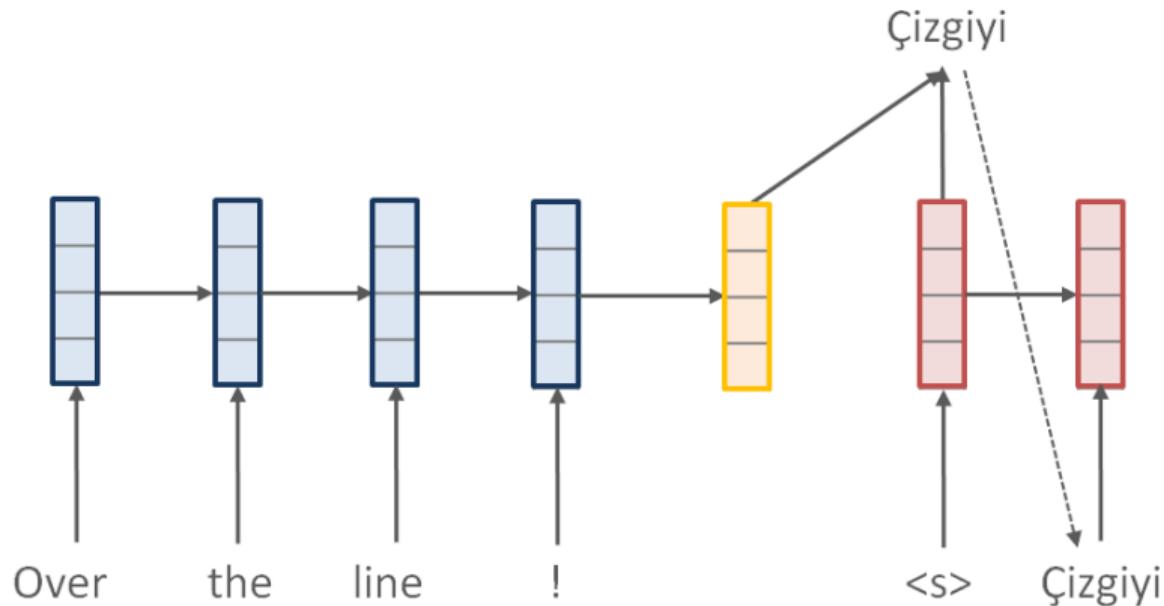
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



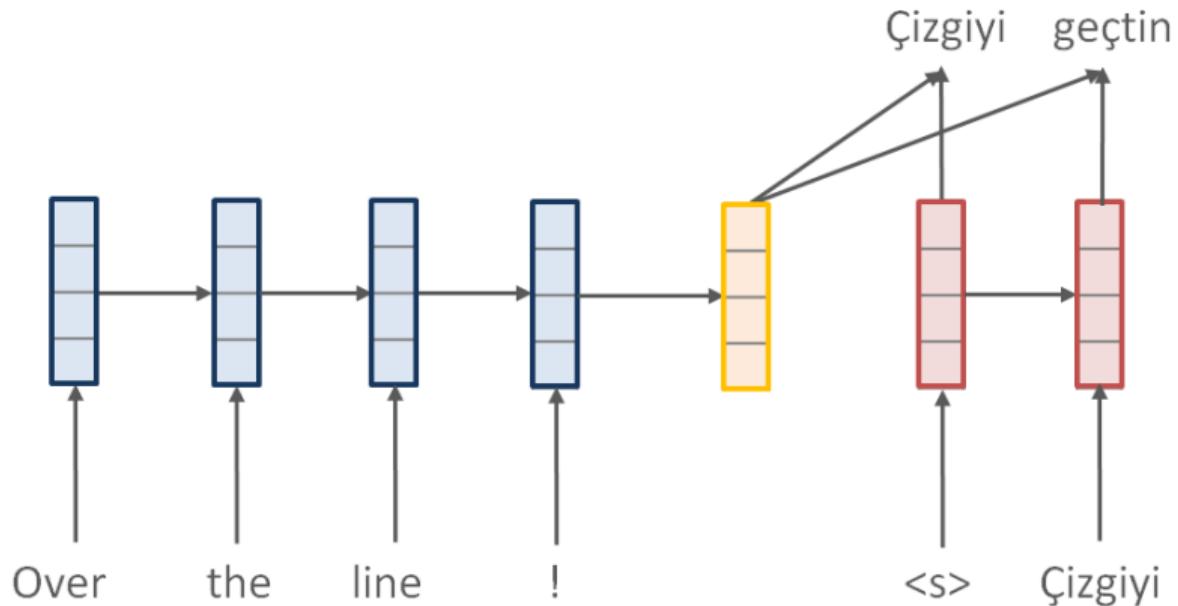
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



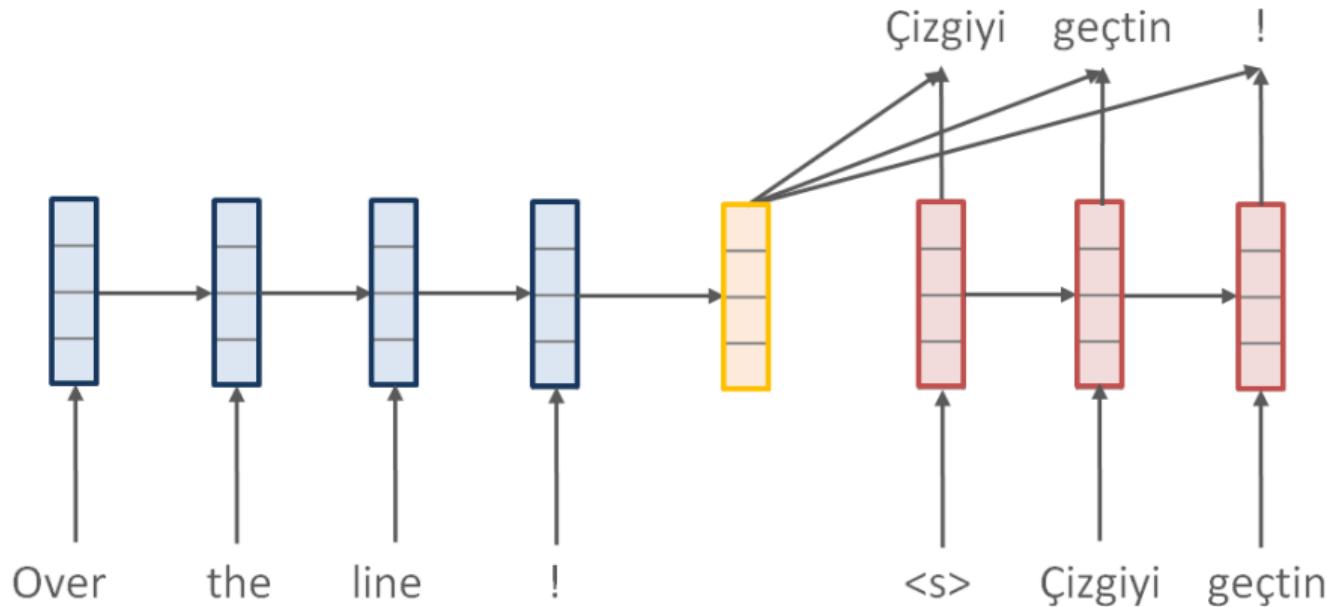
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



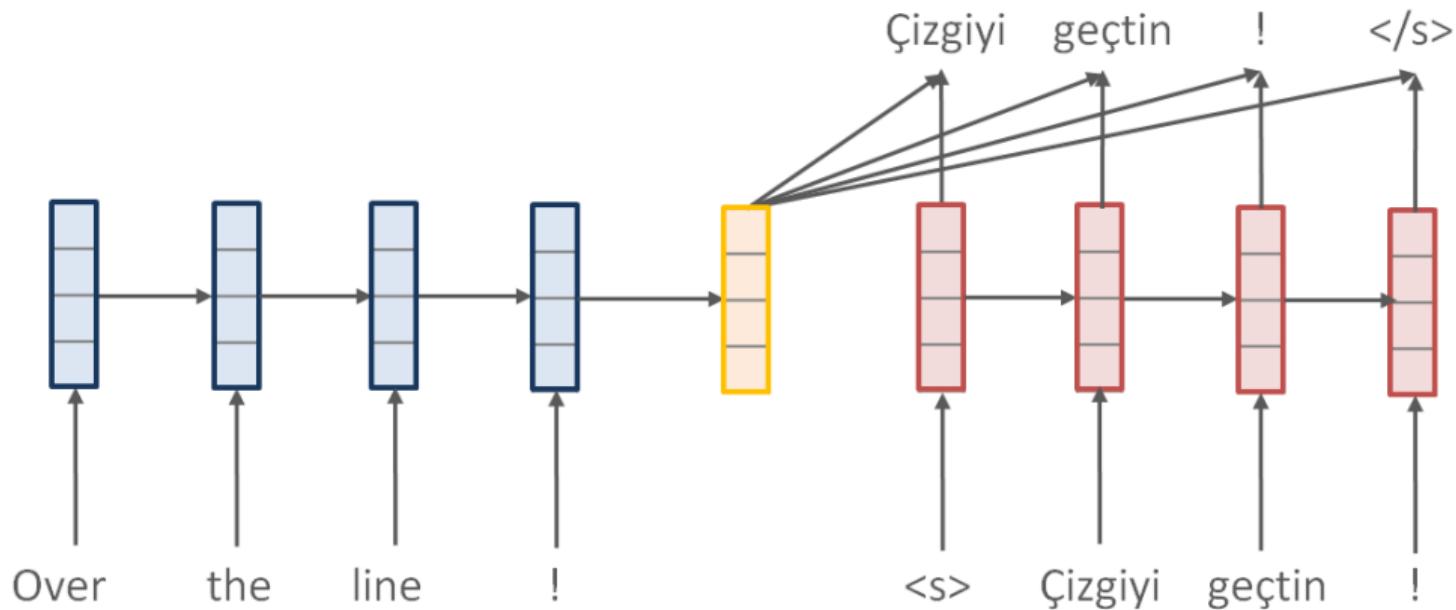
# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



# Recurrent Neural Network 1

$$f(y_{1:T}, x_{1:T}; \theta)$$



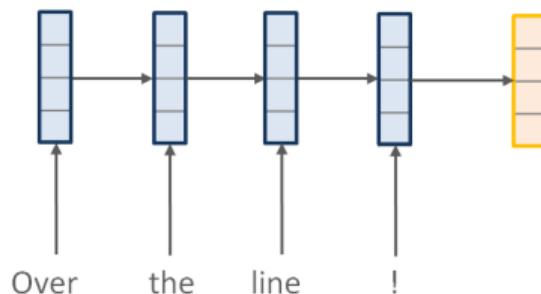
# Recurrent Neural Network Math

Encoder:

$$\mathbf{h}_s^x \leftarrow \text{RNN}(\mathbf{h}_{s-1}^x, x_s)$$

Context:

$$\mathbf{c} = \mathbf{h}_S^x$$



# Recurrent Neural Network Math

Encoder:

$$\mathbf{h}_s^x \leftarrow \text{RNN}(\mathbf{h}_{s-1}^x, x_s)$$

Context:

$$\mathbf{c} = \mathbf{h}_S^x$$

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

Prediction:

$$p(y_{t+1} \mid y_{1:t}, x) = \text{softmax}(\mathbf{W}[\mathbf{h}_t; \mathbf{c}])$$

# Recurrent Neural Network Math

Encoder:

$$\mathbf{h}_s^x \leftarrow \text{RNN}(\mathbf{h}_{s-1}^x, x_s)$$

Context:

$$\mathbf{c} = \mathbf{h}_S^x$$

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

Prediction:

$$p(y_{t+1} \mid y_{1:t}, x) = \text{softmax}(\mathbf{W}[\mathbf{h}_t; \mathbf{c}])$$

Generation Score:

$$f(y_{1:T}, x; \theta) = \log \sum_{t=1}^T p(y_t \mid y_{1:t-1}, x)$$

# What can the decoder learn to say?

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

# What can the decoder learn to say?

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

Toy Example:

Well-balanced parenthesis language with random nesting-level indicators,

- Vocabulary: ( ) 0 1 2 3 4
- Example String: 0 ( ( 2 ) ( ( ( 4 4 4 ) 3 ) ...

Proxy Question: What does  $\mathbf{h}_t$  look like over time?

# LSTMVis - Parenthesis Language (Strobelt et al. [2016] w/ IBM)

Temporary

# What can the decoder learn to say?

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

# What can the decoder learn to say?

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

**Harder Example:** Natural language outputs with complex syntax.

# LSTMVis - Parenthesis Language (Strobelt et al. [2016] w/ IBM)

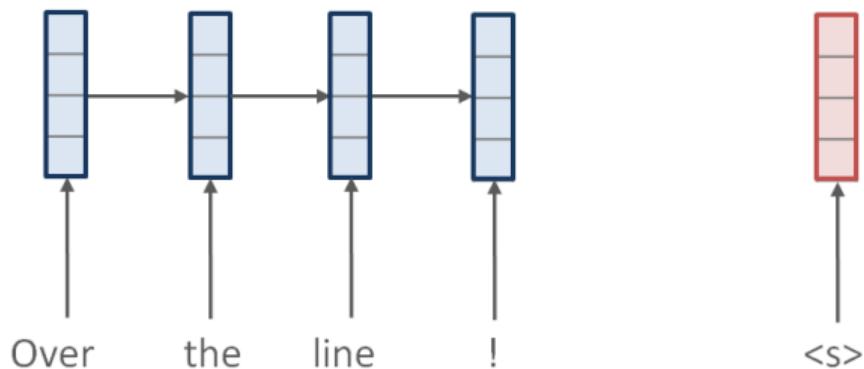
Temporary

# LSTMVis - Natural Language (Strobelt et al. [2016] w/ IBM)

Temporary

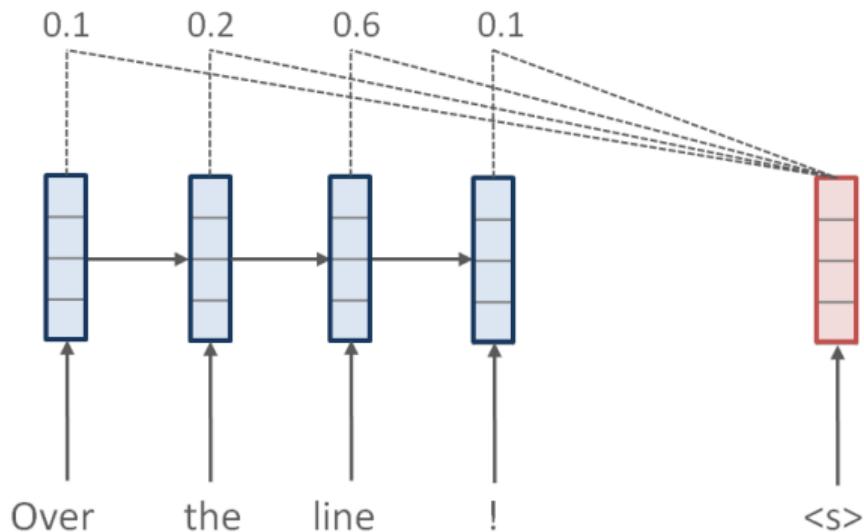
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



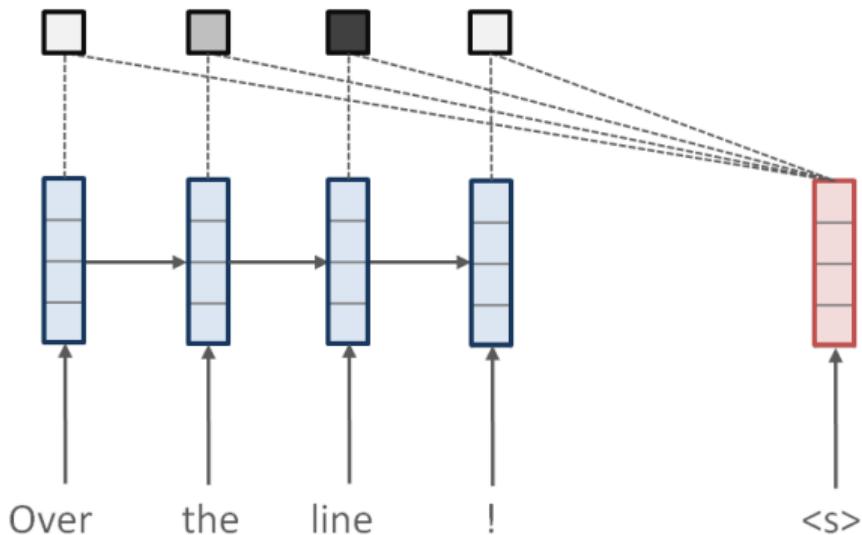
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



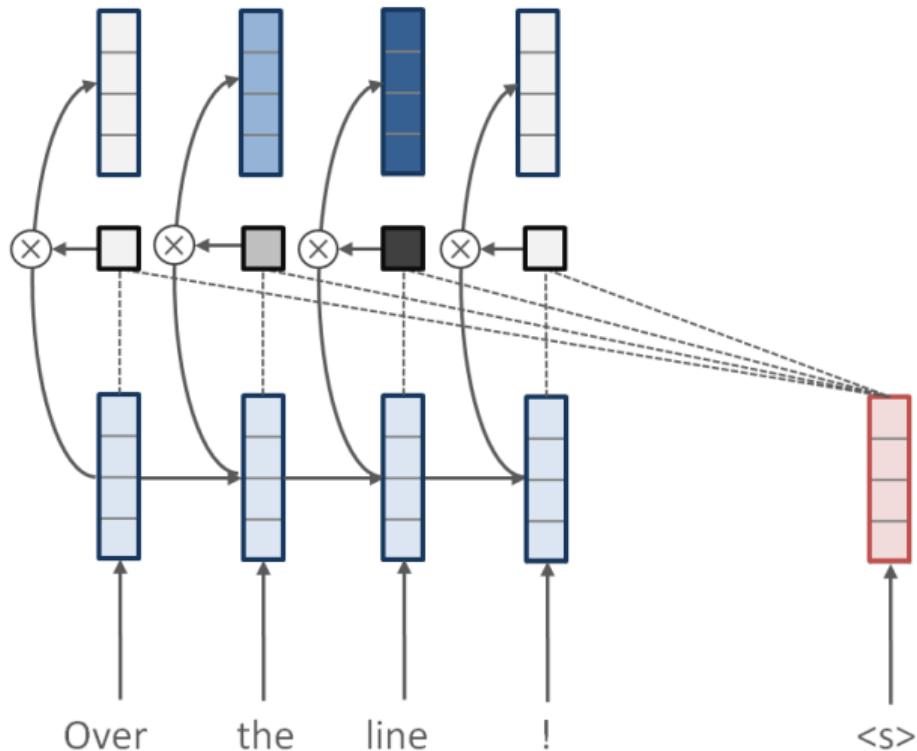
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



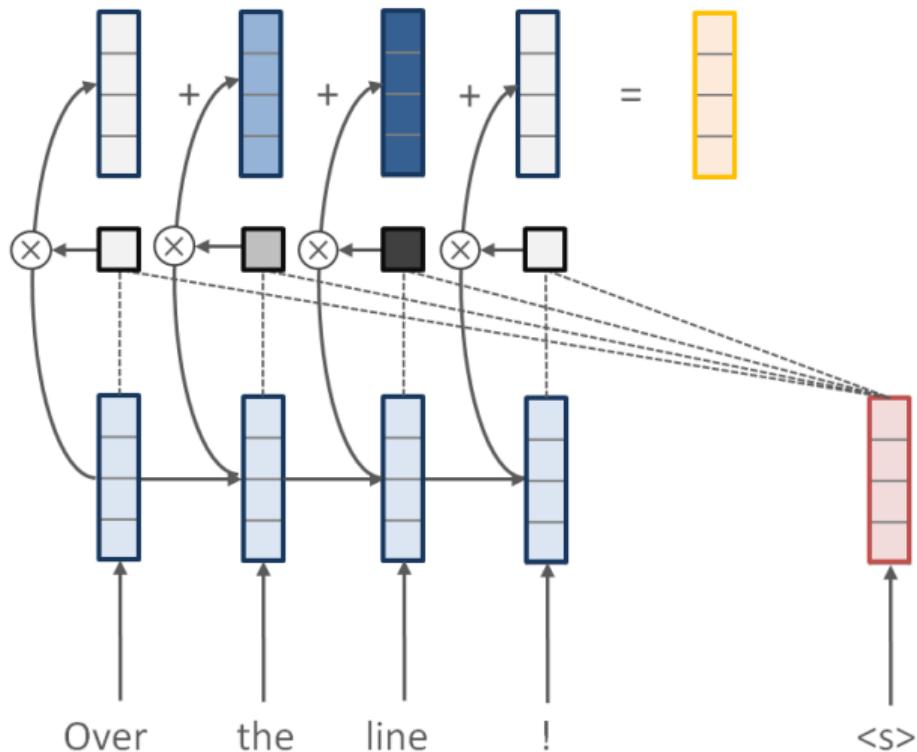
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



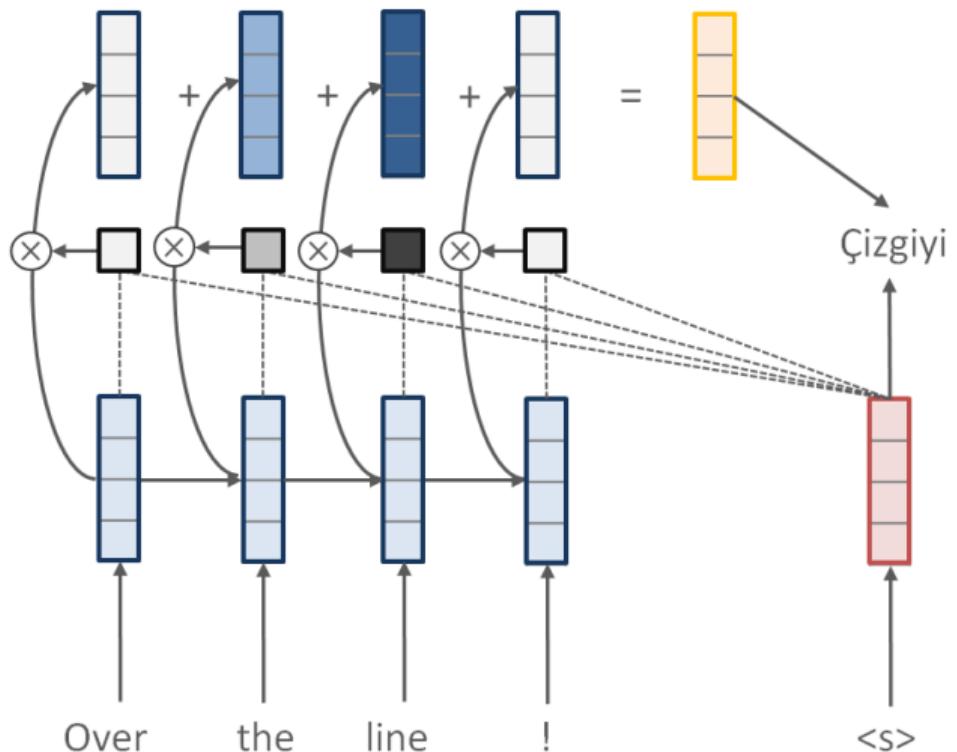
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



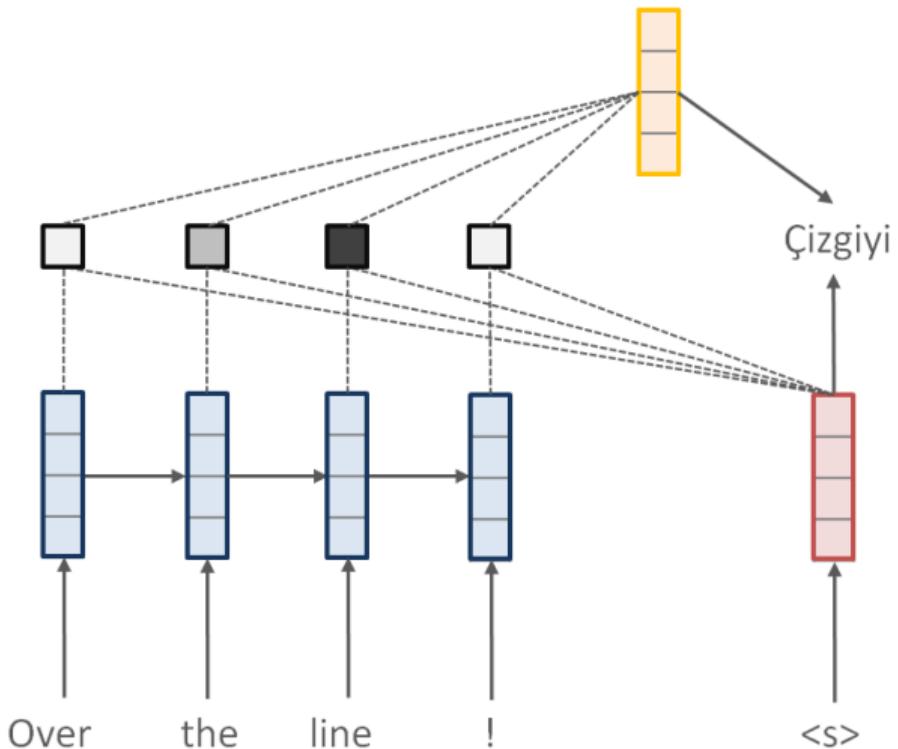
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



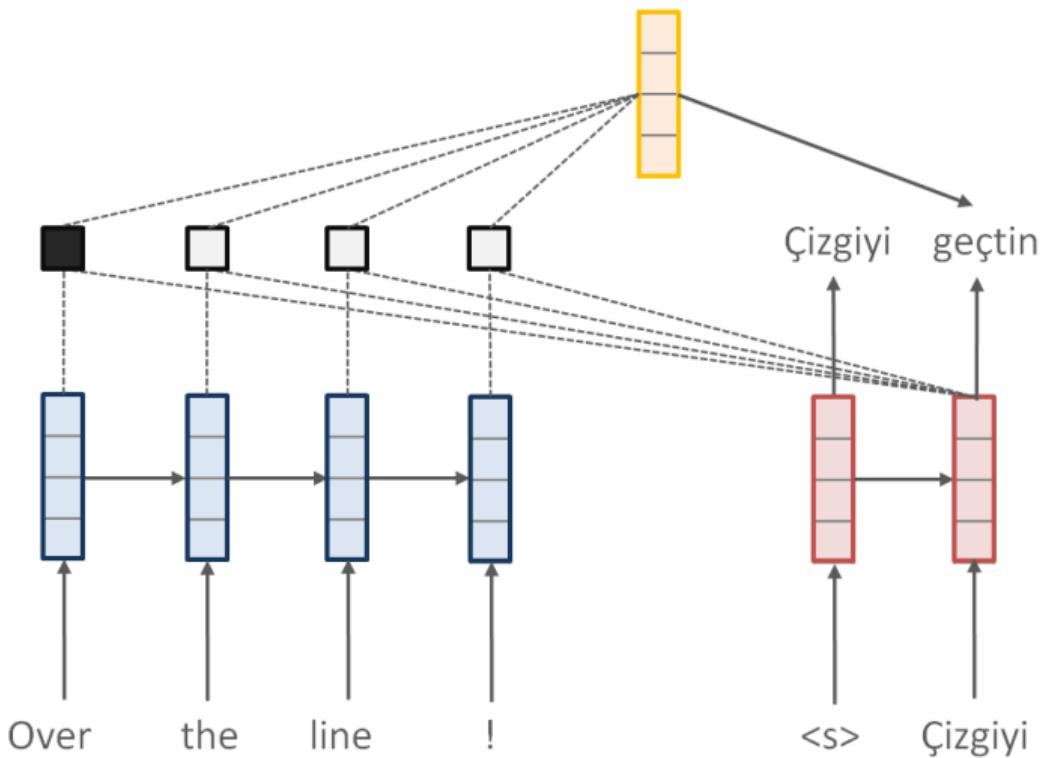
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



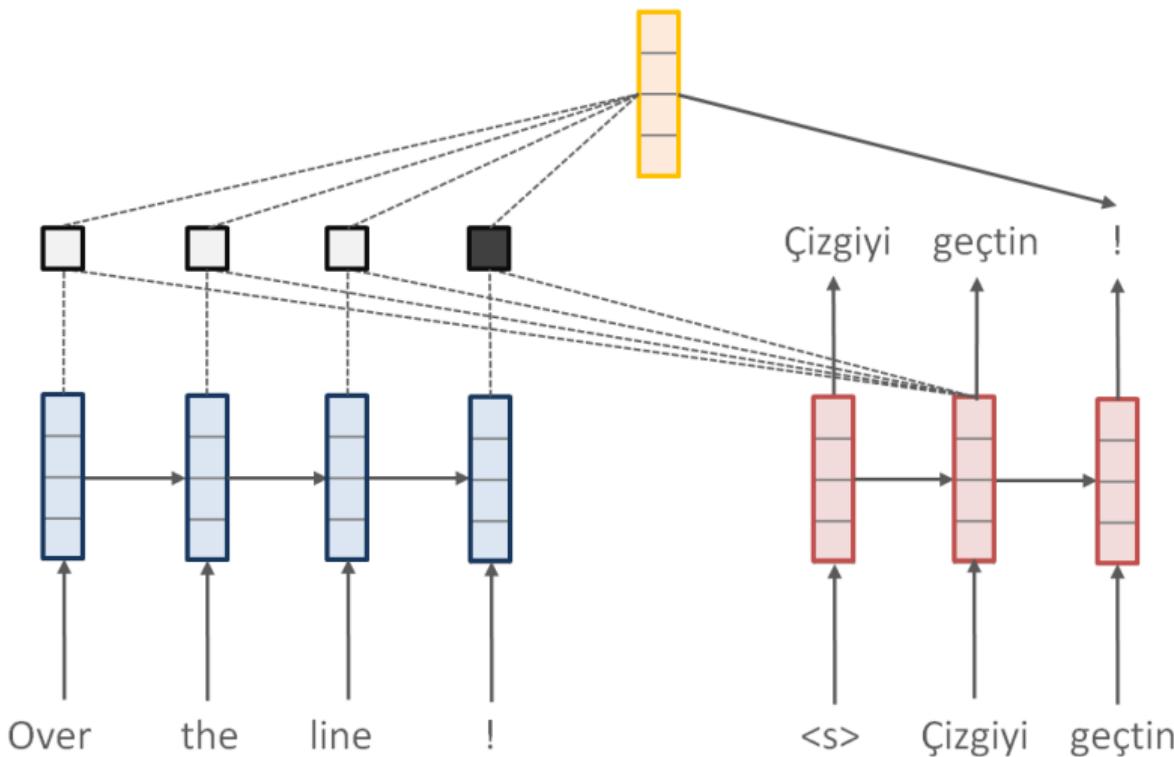
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



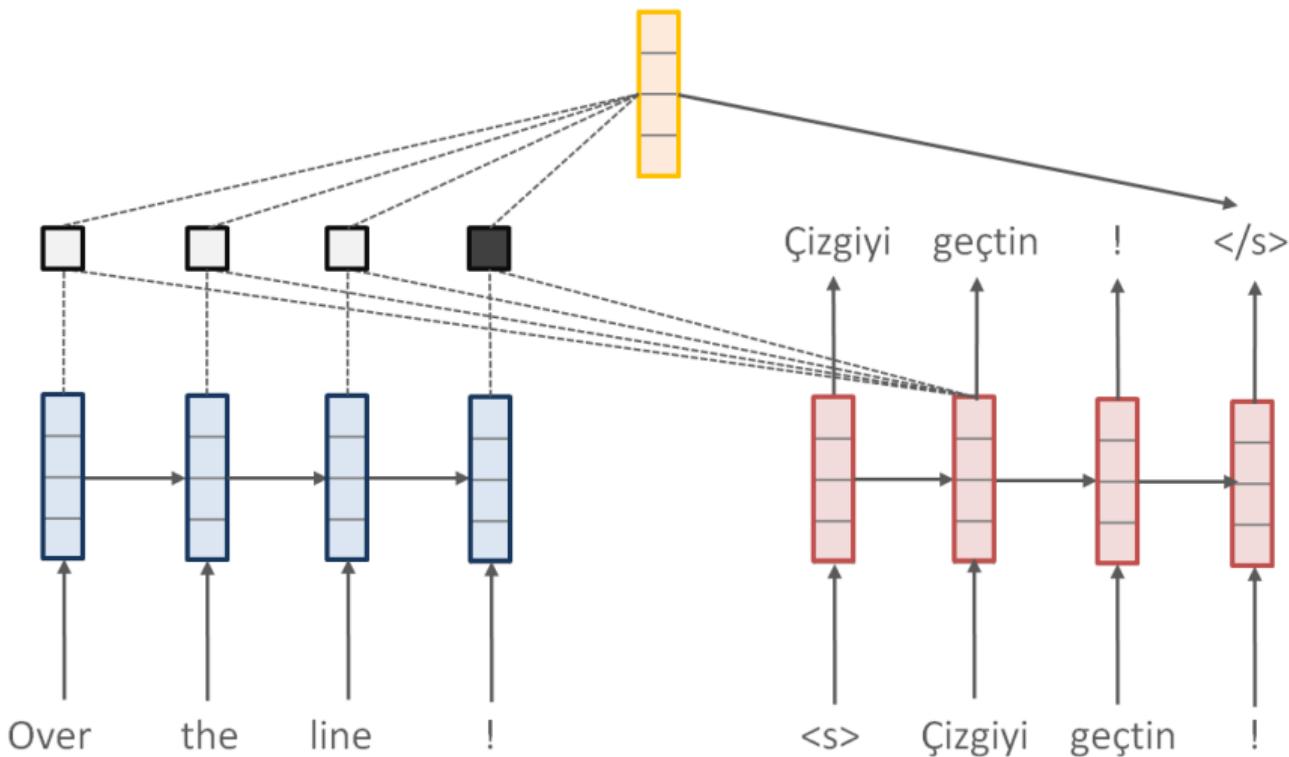
# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



# Recurrent Neural Network 2 - Seq2Seq + Attention

$$f(y_{1:T}, x_{1:T}; \theta)$$



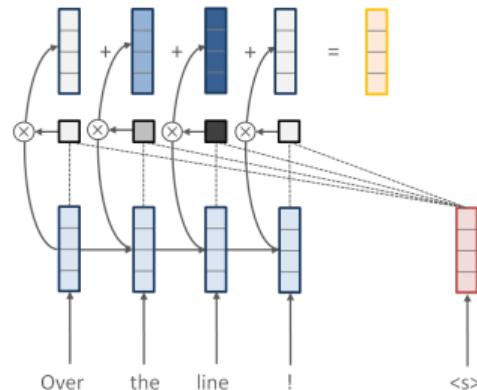
# Attention Math

Encoder:

$$\mathbf{h}_s^x \leftarrow \text{RNN}(\mathbf{h}_{s-1}^x, x_s)$$

Attention (Dynamic Context)

$$\alpha \leftarrow \text{softmax}([\mathbf{h}_1^x; \dots; \mathbf{h}_S^x]^\top \mathbf{h}_t) \quad \mathbf{c} \leftarrow \sum_{s=1}^S \alpha_s \mathbf{h}_s^x$$



# Attention Math

Encoder:

$$\mathbf{h}_s^x \leftarrow \text{RNN}(\mathbf{h}_{s-1}^x, x_s)$$

Attention (Dynamic Context)

$$\alpha \leftarrow \text{softmax}([\mathbf{h}_1^x; \dots; \mathbf{h}_S^x]^\top \mathbf{h}_t) \quad \mathbf{c} \leftarrow \sum_{s=1}^S \alpha_s \mathbf{h}_s^x$$

Decoder:

$$\mathbf{h}_t \leftarrow \text{RNN}(\mathbf{h}_{t-1}, y_t)$$

Prediction:

$$p(y_{t+1} \mid y_{1:t}, x) = \text{softmax}(\mathbf{W}[\mathbf{h}_t; \mathbf{c}])$$

# How does attention control what is said?

## Attention (Dynamic Context)

$$\alpha \leftarrow \text{softmax}([\mathbf{h}_1^x; \dots; \mathbf{h}_S^x]^\top \mathbf{h}_t) \quad \mathbf{c} \leftarrow \sum_{s=1}^S \alpha_s \mathbf{h}_s^x$$

- Can we use this to control the output of the system?
- Can we examine how errors enter into translation?

# Seq2SeqVis

(Strobelt et al. [2019] w/ IBM)

Temporary



An open-source neural machine translation system.

English Français 简体中文 한국어  
日本語 Русский العربية

## Home

[Quickstart \[Lua\]](#)

[Quickstart \[Python\]](#)

[Advanced guide](#)

[Models and Recipes](#)

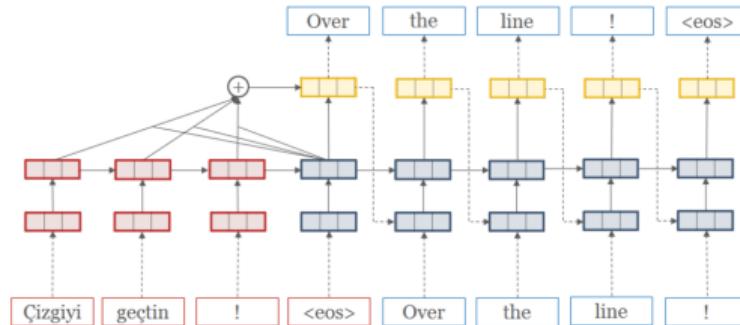
[FAQ](#)

[About](#)

[Documentation](#)

# Home

OpenNMT is a industrial-strength, open-source (MIT) neural machine translation system utilizing the [Torch/PyTorch](#) mathematical toolkit.



OpenNMT is used as provided in [production](#) by major translation providers. The system is designed to be simple to use and easy to extend, while maintaining efficiency and state-of-the-art translation accuracy.



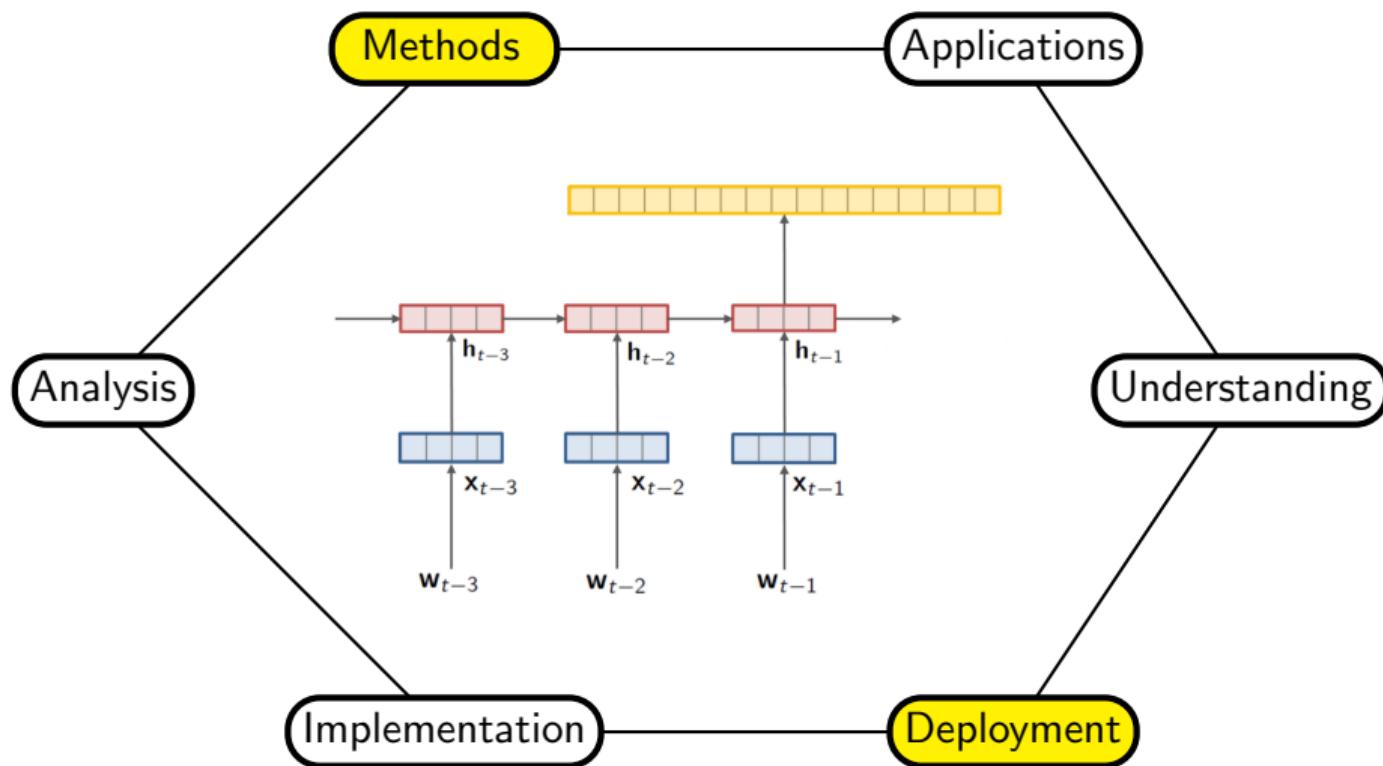
- Collaborative open-source project started at Harvard, now self-sustaining.
- Used in production by Systran, Ubiquis, Booking.com, and others.
- Over 100 developers in France, China, Japan, Portugal, and the US.
- Designed to be research extensible to latest machine translation techniques.
- Pretrained models for translation as well as everything in this talk.

# OpenNMT Workshop

Paris 2018



## Part 3: Structured Modeling



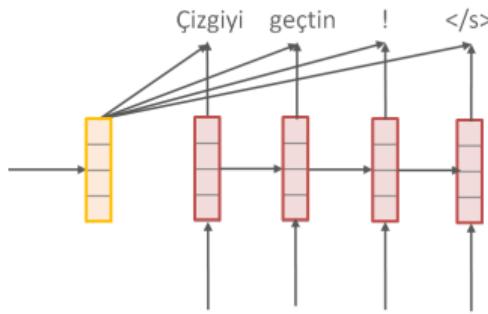
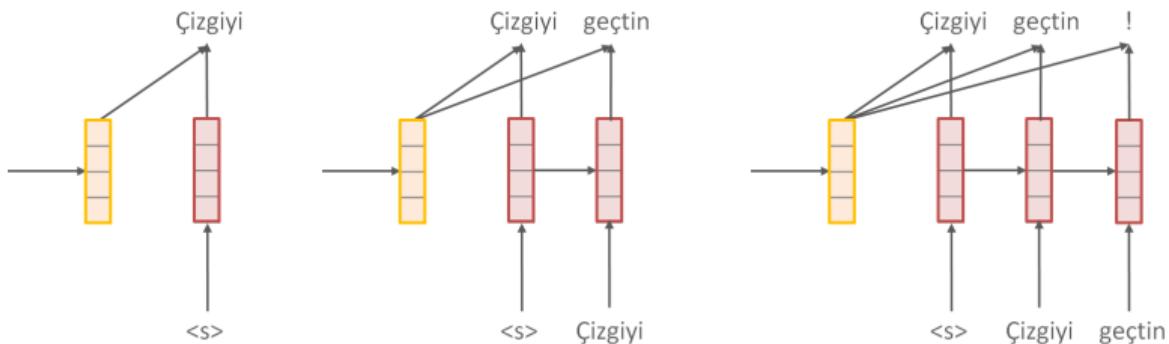
# Machine Learning for Text Generation

$$y_{1:T}^* = \arg \max_{y_{1:T}} f(y_{1:T}; x, \theta)$$

- Input  $x_{1:S}$ , *what to talk about*
- Output text  $y_{1:T}^*$ , *how to say it*
- Model  $f(\cdot; \theta)$ , learned from data

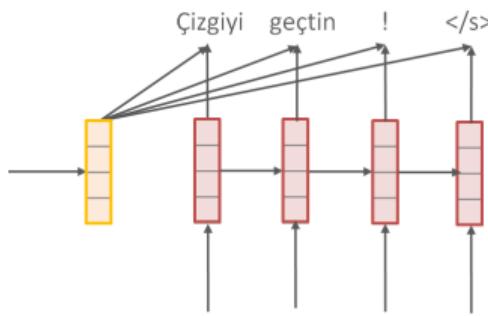
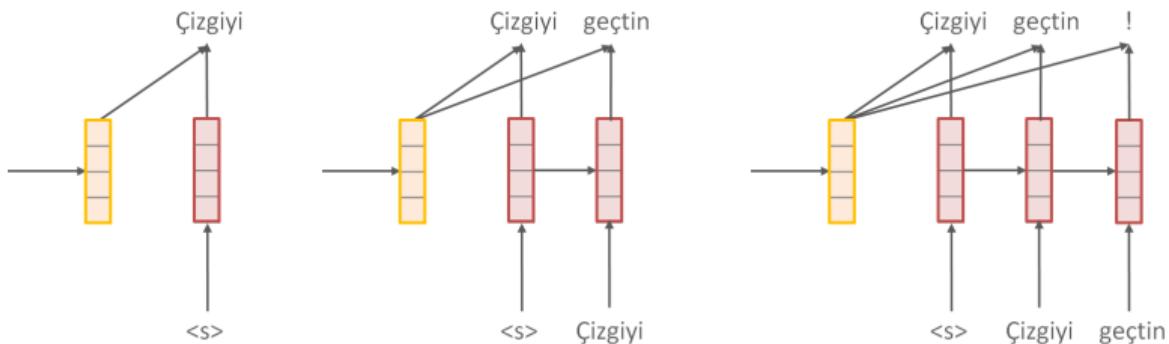
# Training Seq2Seq

Parameters  $\theta$  are trained to predict the next word *given the true history*.



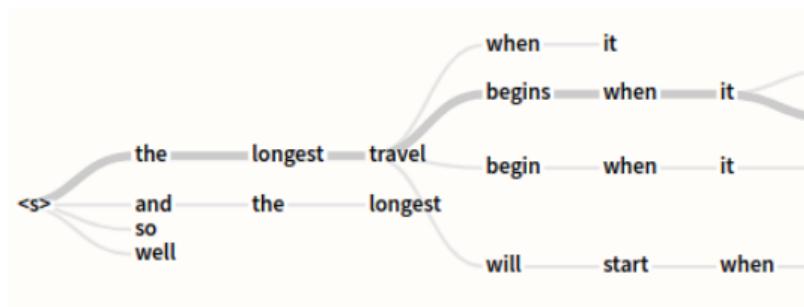
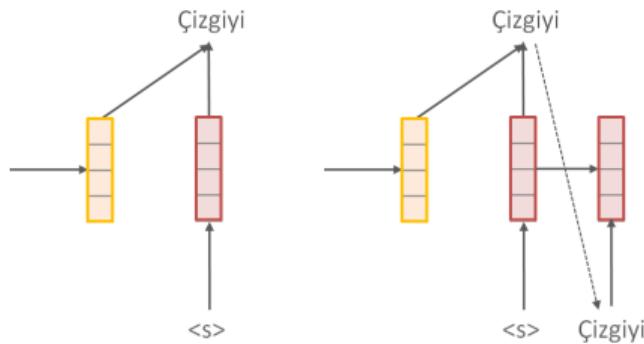
# Training Seq2Seq

Parameters  $\theta$  are trained to predict the next word *given the true history*.



# Deploying Seq2Seq

Parameters  $\theta$  is deployed to predict a next word *given the predicted history*.

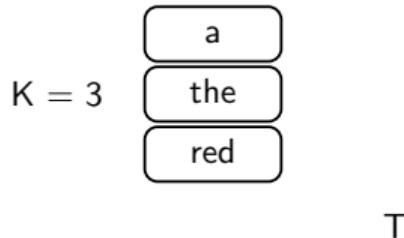


Requires predicting best sequence

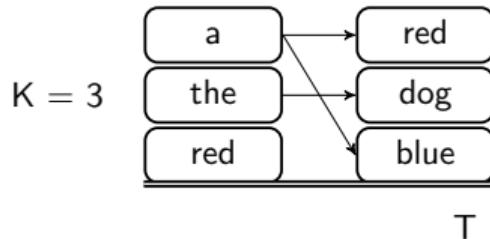
$$y_{1:T}^* = \arg \max_{y_{1:T}} f(y_{1:T}; \theta) = \arg \max_{y_{1:T}} \sum_t \log p(y_t | y_{1:t-1}, \mathbf{c}; \theta)$$

However: Completely intractable for RNNs  $O(\#\text{vocab}^T)$

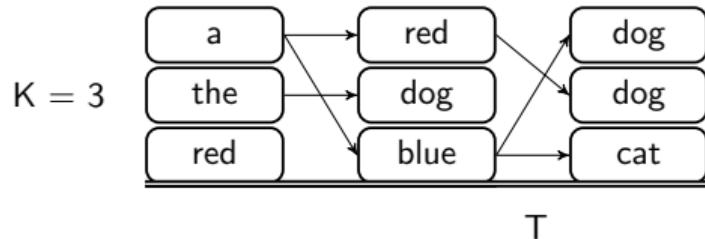
## Standard Heuristic Method: Beam Search



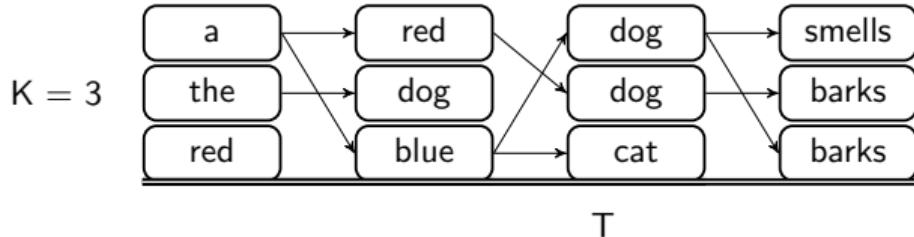
## Standard Heuristic Method: Beam Search



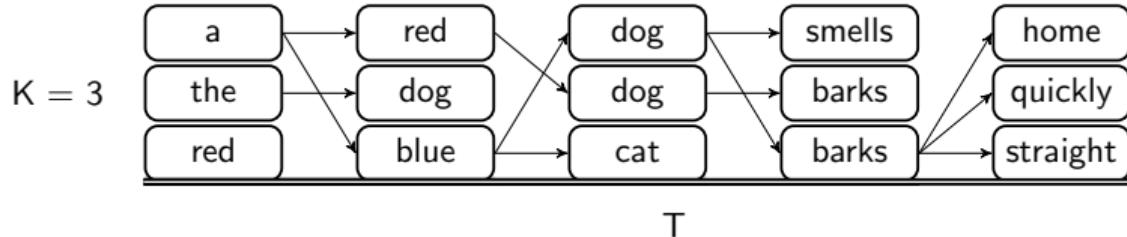
# Standard Heuristic Method: Beam Search



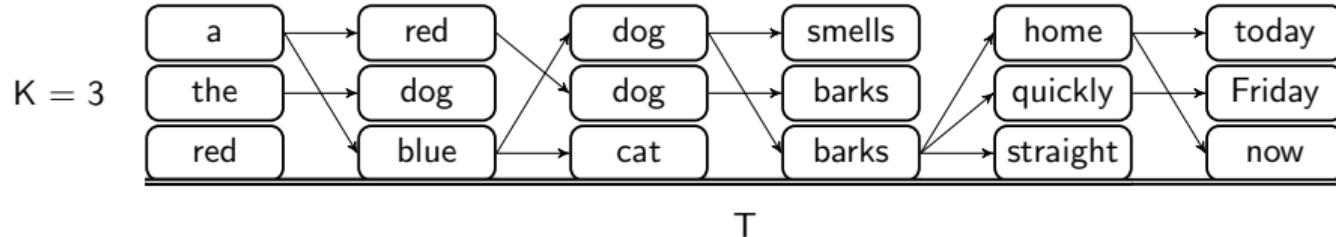
# Standard Heuristic Method: Beam Search



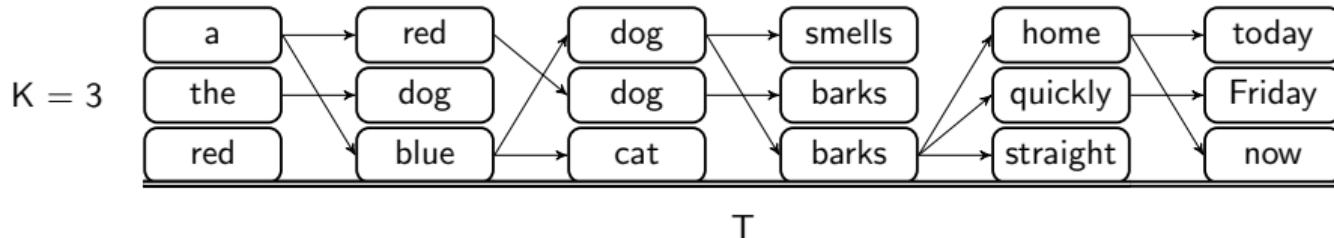
# Standard Heuristic Method: Beam Search



# Standard Heuristic Method: Beam Search



## Standard Heuristic Method: Beam Search



- ① Compute the score of every possible next word.

$$f(y_t, y_{1:t-1}^{(k)}) \leftarrow \log p(y_t \mid y_{1:t-1}^{(k)}, x) + \log p(y_{1:t-1}^{(k)} \mid x)$$

- ② Prune to only the  $K$  highest-scoring,

$$y_{1:t}^{(1:K)} \leftarrow K \arg \max_{y_{1:t}} f(y_t, y_{1:t-1}^{(k)})$$

# Theoretical Issues with RNN-based Models

## ① Label Bias

- Training is locally discriminative, but prediction is over entire sequences.

## ② Exposure Bias

- Training conditions on true history ( $y_{1:t-1}$ ) but generates with predicted history.

# Theoretical Issues with RNN-based Models

## ① Label Bias

- Training is locally discriminative, but prediction is over entire sequences.

## ② Exposure Bias

- Training conditions on true history ( $y_{1:t-1}$ ) but generates with predicted history.

## ③ Metric Bias

- Training uses multiclass classification, but evaluation uses n-gram match.

# Theoretical Issues with RNN-based Models

## ① Label Bias

- Training is locally discriminative, but prediction is over entire sequences.

## ② Exposure Bias

- Training conditions on true history ( $y_{1:t-1}$ ) but generates with predicted history.

## ③ Metric Bias

- Training uses multiclass classification, but evaluation uses n-gram match.

# Research

Can we better model discrete sequences for text generation?

Applications:

- (1) Improvements in training with less supervision.
- (2) Effective methods for downscaling translation models.

# Sequence-to-Sequence Learning as Beam Search Optimization

(Vaswani et al., 2016)

Proposal: Directly modify the RNN training procedure to fix test biases.

- ① Label Bias
- ② Exposure Bias
- ③ Metric Bias

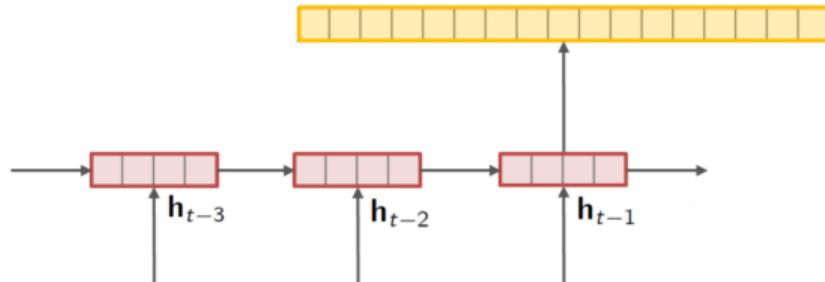
# Modification 1: Global Scoring Function

## Issue: Label Bias

- Training is locally discriminative, but prediction is over entire sequences.

## Proposed Fix:

- Replace  $\log p(y_t|y_{1:t-1}^{(k)}, \mathbf{c}; \theta)$  with a directly learned function  $f(y_t, y_{1:t-1}^{(k)}, x; \theta)$



## Modification 2: Beam Search at Training

**Issue:** Exposure Bias

- Training conditions on true history  $(y_{1:t-1})$  but generates with predicted history.

**Proposed Fix:** During training:

- ① Compute the score of every possible next word.

$$f(y_t, y_{1:t-1}^{(k)}) \leftarrow \log p(y_t \mid y_{1:t-1}^{(k)}, \mathbf{c}) + \log p(y_{1:t-1}^{(k)} \mid \mathbf{c})$$

- ② Prune to only the  $K$  highest-scoring,

$$y_{1:t}^{(1:K)} \leftarrow K \arg \max_{y_{1:t}} f(y_t, y_{1:t-1}^{(k)})$$

## Modification 2: Beam Search at Training

**Issue:** Exposure Bias

- Training conditions on true history  $(y_{1:t-1})$  but generates with predicted history.

**Proposed Fix:** During training:

- ① Compute the score of every possible next word.

$$f(y_t, y_{1:t-1}^{(k)}) \leftarrow f(y_t, y_{1:t-1}^{(k)}, \mathbf{c}; \theta)$$

- ② Prune to only the  $K$  highest-scoring,

$$y_{1:t}^{(1:K)} \leftarrow K \arg \max_{y_{1:t}} f(y_t, y_{1:t-1}^{(k)})$$

## Modification 3: Train with Margin

**Issue:** Metric Bias

- Training uses multiclass classification, but evaluation uses n-gram match.

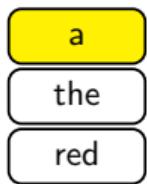
**Proposed Fix:** Use a structured SVM-style training loss:

- Margin between ground truth sequence  $\hat{y}$  and worst predicted sequence  $y^{(K)}$

$$\mathcal{L}(\theta) = \sum_t \Delta(\hat{y}_{1:t}, y_{1:t}^K) \left[ 1 - f(\hat{y}_t, y_{1:t-1}^{(g)}, \mathbf{c}) + f(y_t^{(K)}, y_{1:t-1}^{(K)}, \mathbf{c}) \right]$$

- Slack-rescaled, margin-based sequence criterion, at each time step.
- $\Delta$  is a task specific sequence cost, i.e. ngram-mismatch

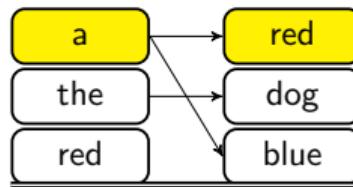
## Beam Search Optimization Example



- Color **True**: ground-truth sequence  $\hat{y}$
- Color **Red**: last sentence  $y^{(K)}$  upon violation

Strategy: upon violation, restart from ground truth (learning as search optimization ?)

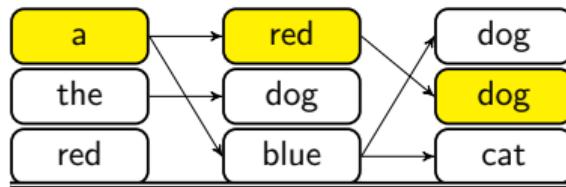
# Beam Search Optimization Example



- Color **True**: ground-truth sequence  $\hat{y}$
- Color **Red**: last sentence  $y^{(K)}$  upon violation

Strategy: upon violation, restart from ground truth (learning as search optimization ?)

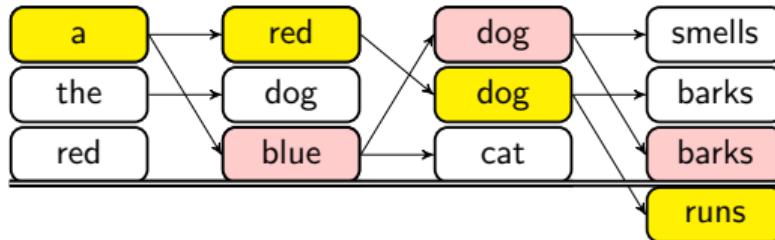
# Beam Search Optimization Example



- Color **True**: ground-truth sequence  $\hat{y}$
- Color **Red**: last sentence  $y^{(K)}$  upon violation

Strategy: upon violation, restart from ground truth (learning as search optimization ?)

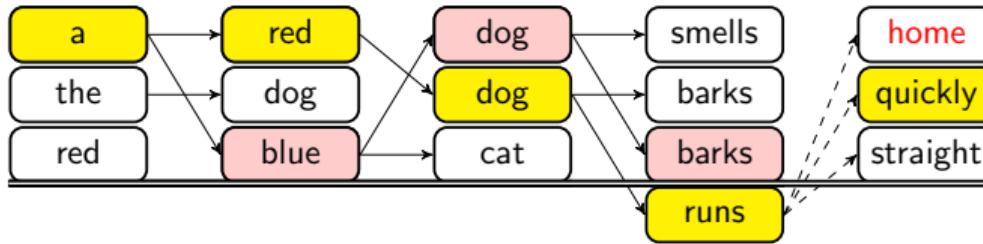
# Beam Search Optimization Example



- Color **True**: ground-truth sequence  $\hat{y}$
- Color **Red**: last sentence  $y^{(K)}$  upon violation

Strategy: upon violation, restart from ground truth (learning as search optimization ?)

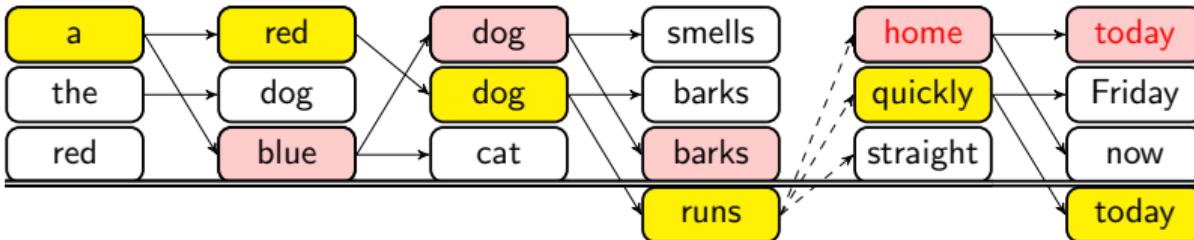
# Beam Search Optimization Example



- Color **True**: ground-truth sequence  $\hat{y}$
- Color **Red**: last sentence  $y^{(K)}$  upon violation

Strategy: upon violation, restart from ground truth (learning as search optimization ?)

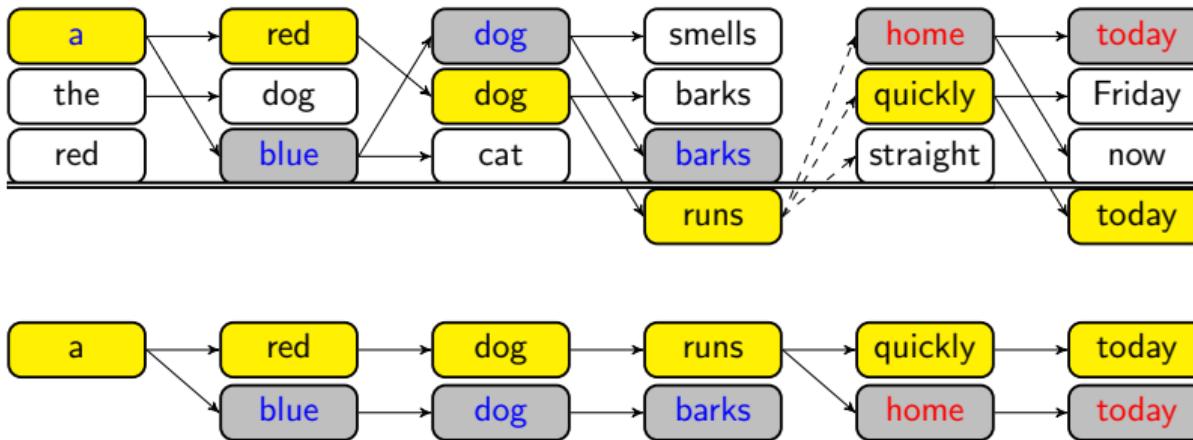
# Beam Search Optimization Example



- Color **True**: ground-truth sequence  $\hat{y}$
- Color **Red**: last sentence  $y^{(K)}$  upon violation

Strategy: upon violation, restart from ground truth (learning as search optimization ?)

# Parameter Updates: Structured Backpropagation



- Margin gradients are sparse, only violating sequences get updates.
- Backprop as efficient as standard models, avoid exponential sum.

# Results

Train Beam	$K = 1$	$K = 5$	$K = 10$
Word Ordering (BLEU)			
seq2seq	25.2	29.8	31.0
BSO	28.0	33.2	34.3
BSO-Con	<b>28.6</b>	<b>34.3</b>	<b>34.5</b>

# Results

Train Beam	$K = 1$	$K = 5$	$K = 10$
Word Ordering (BLEU)			
seq2seq	25.2	29.8	31.0
BSO	28.0	33.2	34.3
BSO-Con	<b>28.6</b>	<b>34.3</b>	<b>34.5</b>

# Results

Train Beam	$K = 1$	$K = 5$	$K = 10$
Word Ordering (BLEU)			
seq2seq	25.2	29.8	31.0
BSO	28.0	33.2	34.3
BSO-Con	<b>28.6</b>	<b>34.3</b>	<b>34.5</b>
Machine Translation (BLEU)			
seq2seq	22.53	24.03	23.87
BSO, SB- $\Delta$	<b>23.83</b>	<b>26.36</b>	<b>25.48</b>
XENT	17.74	$\leq 20.5$	$\leq 20.5$
DAD	20.12	$\leq 22.5$	$\leq 23.0$
MIXER	20.73	-	$\leq 22.0$

# Sequence Knowledge Distillation (Kim and Rush [2016])

Proposal: Shrink the size of text generation models.

Goal: Replicate knowledge distillation results from multiclass image recognition.

- **Knowledge Distillation:** Train a *student* model to learn from a *teacher* model ???.



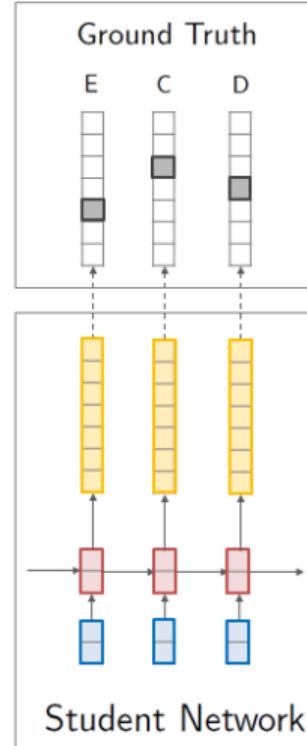
# Baseline

Minimize :

$$\mathcal{L}(\theta) = -\sum_t \log p(y_t = \hat{y}_t \mid \hat{y}_{1:t-1}, x; \theta)$$

where  $\hat{y}_t$  is the ground truth word at time  $t$ .

Cross-entropy with ground truth.

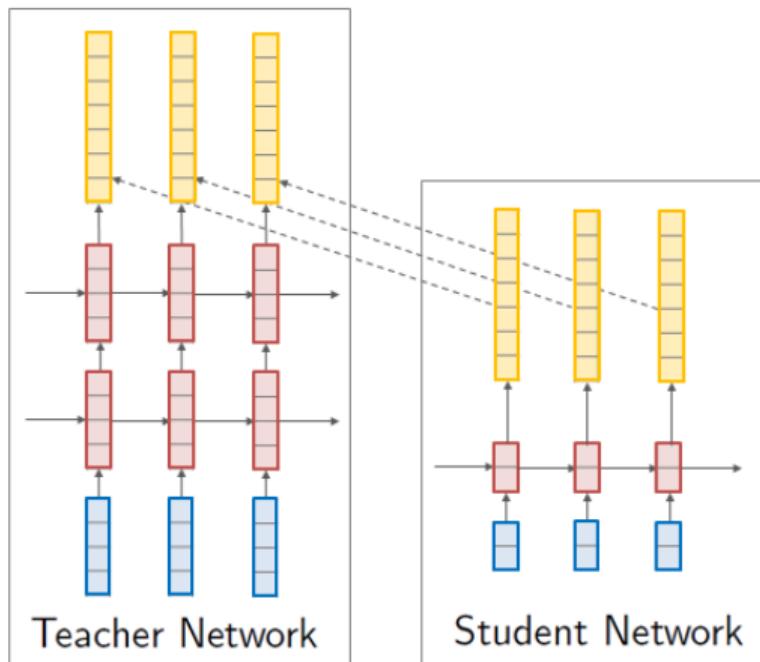


# Multiclass Style: Word-Level Knowledge Distillation

Teacher model:  $q(y_t | y_{1:t-1}, x; \theta_T)$

Cross-entropy between teacher and student

$$\mathcal{L}_{\text{WORD-KD}}(\theta) = - \sum_t \sum_v q(y_t = v | \hat{y}_{1:t-1}, x; \theta_T) \times \log p(y_t = v | \hat{y}_{1:t-1}, x; \theta)$$



# Sequence-Level Knowledge Distillation

**Motivation:** Replace multiclass with sequence-level cross-entropy.

$$\mathcal{L}_{\text{WORD-KD}}(\theta) = - \sum_t \sum_v q(y_t = v \mid \hat{y}_{1:t-1}, x; \theta_T) \times \log p(y_t = v \mid \hat{y}_{1:t-1}, \mathbf{c}; \theta)$$



$$\mathcal{L}_{\text{SEQ-KD}}(\theta) = - \sum_{v_1} \dots \sum_{v_T} q(y_{1:T} = v_{1:T} \mid x; \theta_T) \times \log p(y_{1:T} = v_{1:T} \mid x; \theta)$$

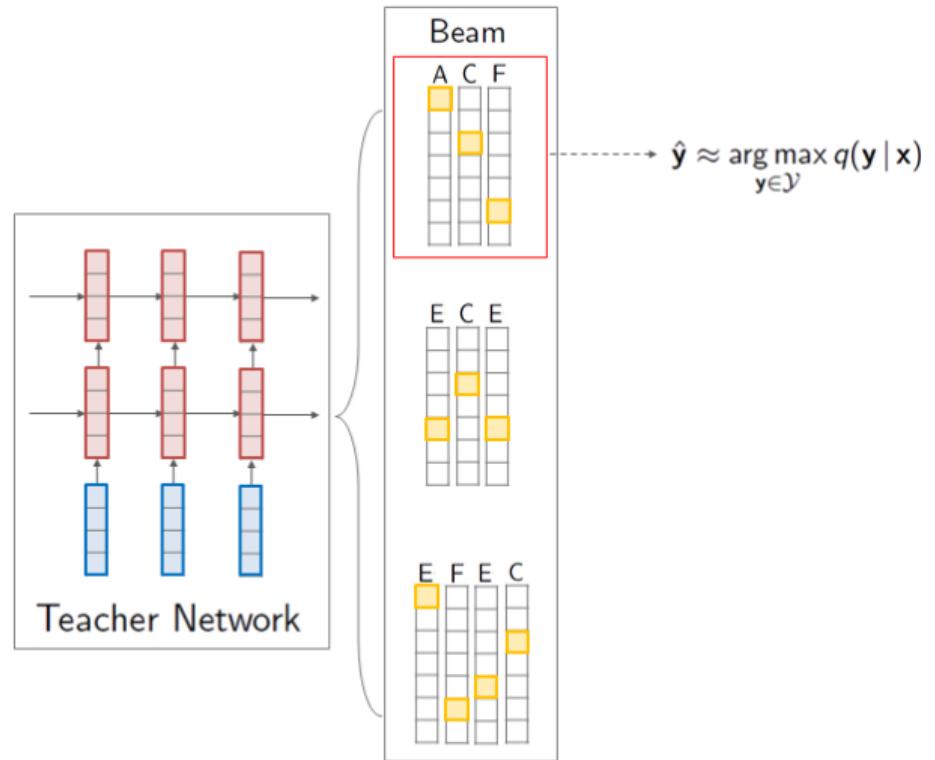
Mimic sequence output of teacher model.

Note: bottom distribution is again intractable.

# Sequence Heuristic Approximation

Approximate  $q(y_{1:T} | x)$  with (beam search) mode sample

$$q(y_{1:T} | x) \approx \mathbf{1}_{\{y \in \mathcal{Y}\}} \{\arg \max_y q(y_{1:T} | x)\}$$

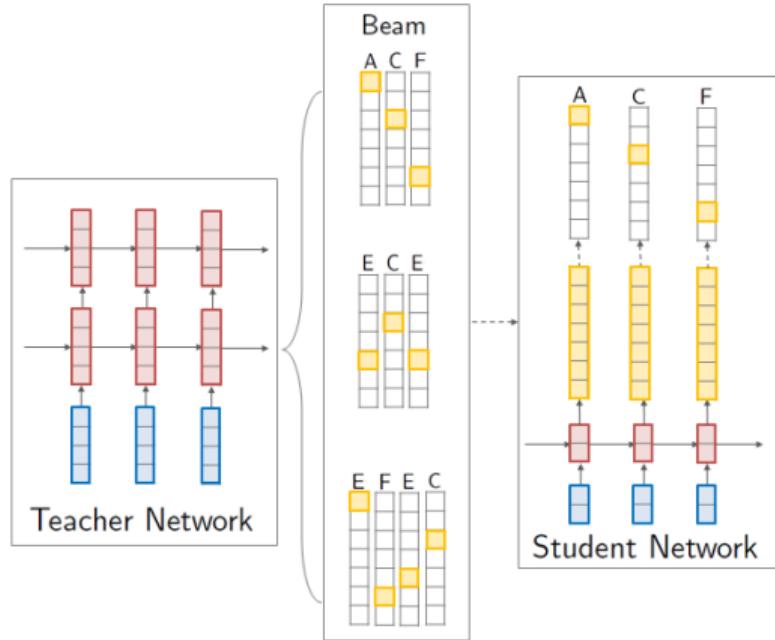


# Sequence-Level Knowledge Distillation

$$\mathcal{L}_{\text{SEQ-KD}}(\theta) = -\log p(y_{1:T}^* | x; \theta)$$

$$\approx - \sum_{v_{1:T}} q(y_{1:T} = v_{1:T} | x; \theta_T) \log p(y_{1:T} | x; \theta)$$

Extension:  $\mathcal{L}_{\text{SEQ-INTER}}(\theta)$  select sample based on ground truth  $\hat{y}$  as well.



## Results: WMT English → German Translation

Model	$\text{BLEU}_{K=1}$	$\Delta_{K=1}$	$\text{BLEU}_{K=5}$	$\Delta_{K=5}$
$4 \times 1000$				
Teacher	17.7	—	19.5	—

## Results: WMT English → German Translation

Model	$\text{BLEU}_{K=1}$	$\Delta_{K=1}$	$\text{BLEU}_{K=5}$	$\Delta_{K=5}$
<hr/>				
$4 \times 1000$				
<hr/>				
Teacher	17.7	—	19.5	—
<hr/>				
$2 \times 500$				
<hr/>				
Student	14.7	—	17.6	—
Word-KD	15.4	+0.7	17.7	+0.1

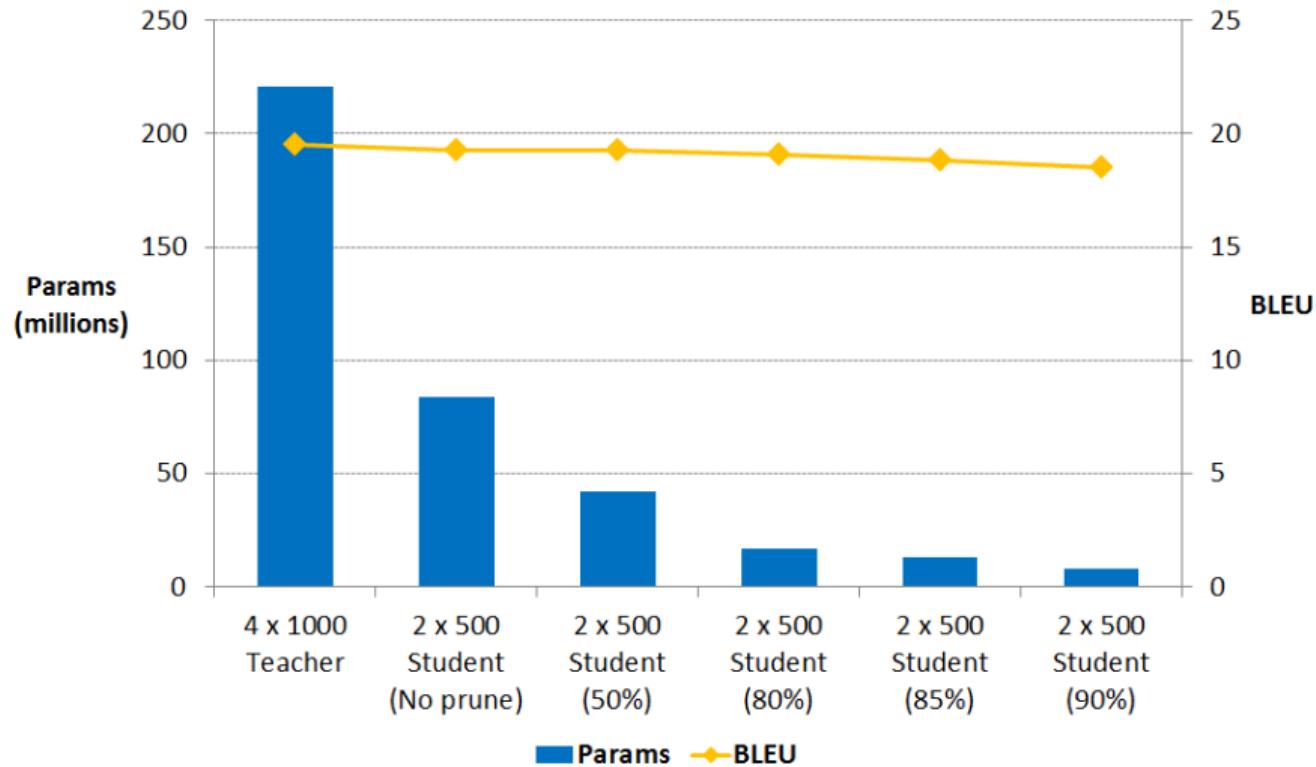
## Results: WMT English → German Translation

Model	BLEU <sub>K=1</sub>	$\Delta_{K=1}$	BLEU <sub>K=5</sub>	$\Delta_{K=5}$
4 × 1000				
Teacher	17.7	—	19.5	—
2 × 500				
Student	14.7	—	17.6	—
Word-KD	15.4	+0.7	17.7	+0.1
Seq-KD	18.9	<b>+4.2</b>	19.0	+1.4
Seq-Inter	18.9	<b>+4.2</b>	19.3	<b>+1.7</b>

## Results: WMT English → German Translation

Model	BLEU <sub>K=1</sub>	$\Delta_{K=1}$	BLEU <sub>K=5</sub>	$\Delta_{K=5}$
4 × 1000				
Teacher	17.7	—	19.5	—
2 × 500				
Student	14.7	—	17.6	—
Word-KD	15.4	+0.7	17.7	+0.1
Seq-KD	18.9	<b>+4.2</b>	19.0	+1.4
Seq-Inter	18.9	<b>+4.2</b>	19.3	<b>+1.7</b>
4 × 1000				
Seq-Inter	19.6	+1.9	19.8	+0.3

# Combining Knowledge Distillation and Pruning

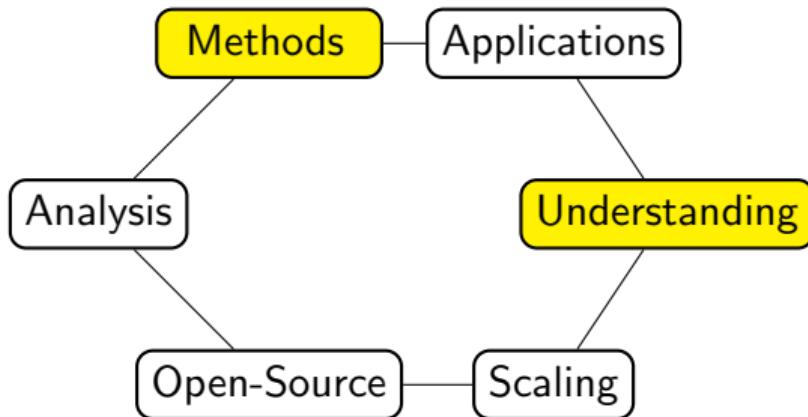


# Application

# WNMT Translation Scaling Shared Task 2018

(Results)

## Part 4: Deep Latent-Variable Mdoels



# Deep Latent-Variable Models

Goal: Extend text generation to Expose specific choices as *discrete* latent variables.

$$p(y, z|x; \theta).$$

# Deep Latent-Variable Models

Goal: Extend text generation to Expose specific choices as *discrete* latent variables.

$$p(y, z|x; \theta).$$

- $y$  is our text output sequence
- $z$  is a collection of latent variables
- $\theta$  are the neural network parameters.

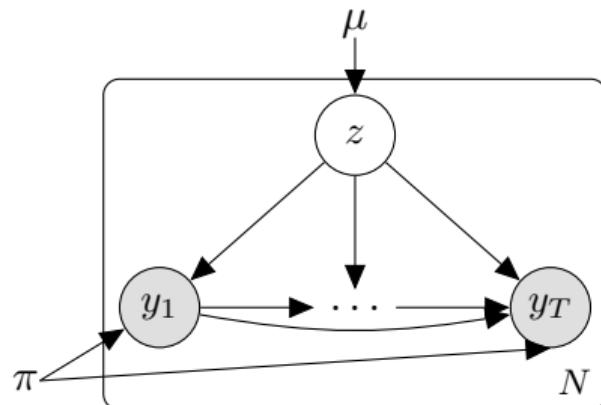
# Example Model: Mixture of RNNs

Generative process:

- ① Draw cluster  $z \in \{1, \dots, K\}$  from a Categorical.
- ② Draw words  $y_{1:T}$  from RNNLM with parameters  $\pi_z$ .

$$p(y, z|x; \theta) = \mu_z \times \text{RNNLM}(y_{1:T}; \pi_z)$$

j



## Posterior Inference

We'll be interested in the *posterior* over latent variables  $z$ :

$$p(z \mid y, x; \theta) = \frac{p(y, z \mid x; \theta)}{p(y \mid x; \theta)} = \frac{p(y \mid x, z; \theta)p(z \mid x; \theta)}{\sum_{z'} p(y \mid x, z'; \theta)p(z' \mid x; \theta)}.$$

## Posterior Inference

We'll be interested in the *posterior* over latent variables  $z$ :

$$p(z | y, x; \theta) = \frac{p(y, z | x; \theta)}{p(y | x; \theta)} = \frac{p(y | x, z; \theta)p(z | x; \theta)}{\sum_{z'} p(y | x, z'; \theta)p(z' | x; \theta)}.$$

How?

- Sum out over all discrete choices (e.g. run  $K$  RNNs).
- Variational inference based methods.

# Application: Summary with Copy-Attention

(Gu et al, 2016) (Gulcehre et al, 2016)

Let  $z$  be a binary latent variable.

- If  $z = 1$ , let the model generate a new word.
- If  $z = 0$ , let the model copy a word from the source.

Inference:

## Pointer-generator model + coverage summary

francis saili has signed a two-year deal to join munster later this year .  
the 24-year-old was part of the new zealand under-20 side that won the junior world championship in italy in 2011 .  
saili 's signature is something of a coup for munster and head coach anthony foley .

(See et al, 2017)

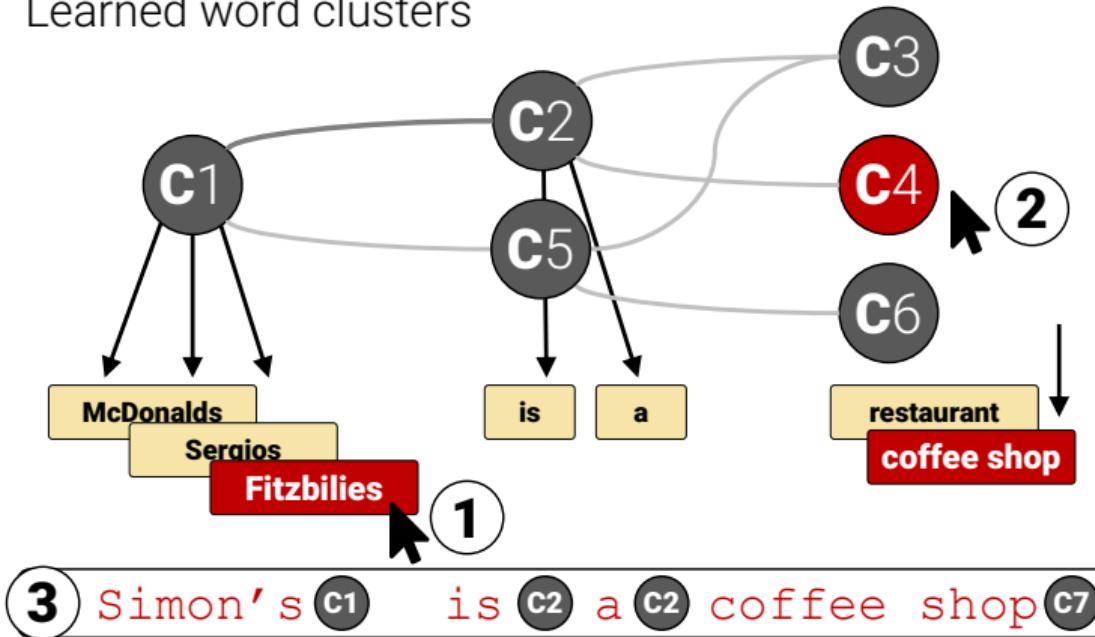
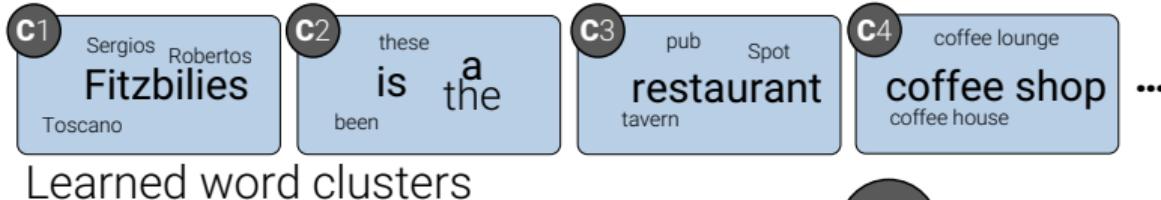
# Latent Variable Models for Generation

**Ongoing Work:** Can we develop other discrete latent-variable models for generation?

**Goals:**

- Model Control
- Model Debugging
- Model Uncertainty

## Example: Learning Neural Templates for Generation



---

**MR** name[The Golden Palace], eatType[coffee shop], food[Fast food],  
priceRange[cheap], customer rating[5 out of 5], area[riverside]

---

**Reference** A coffee shop located on the riverside called The Golden Palace,  
has a 5 out of 5 customer rating. Its price range are fairly cheap  
for its excellent Fast food.

---

# Standard Approach

## Step 1: Encode the Source

Fitzbillies,type[coffee shop],price[< £20],food[Chinese],rate[3/5],area[city centre]

## Step 2: Generate with RNN Decoder

Fitzbillies is a coffee shop providing Chinese food in the moderate price range . It is located in the city centre . Its customer rating is 3 out of 5.

## Issues

- ① Interpretable in its content selection?

*Decisions may come from anywhere in the source  $x$ .*

- ② Controllable in terms of style and form?

*Rely on a learned system to determine content.*

# Neural Template Generation Approach

## Step 1: Encode the Source

Fitzbillies,ty[coffee shop],pr[< £20],food[Chinese],cust[3/5],area[city centre]

# Neural Template Generation Approach

## Step 1: Encode the Source

Fitzbillies,ty[coffee shop],pr[< £20],food[Chinese],cust[3/5],area[city centre]

## Step 2: Select a Template

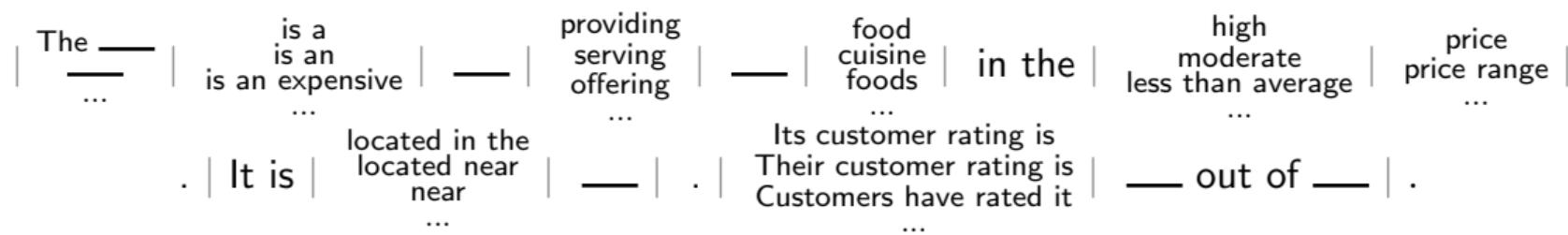
The — | is a — | providing — | food — | high — | price —  
— | is an — | serving — | cuisine — | moderate — | price range —  
... | ... | offering — | foods — | less than average — | ...  
| ... | ... | ... | ... | ... | ...  
. | It is | located in the — | Its customer rating is — | . | . | . |  
| located near — | Their customer rating is — | . |  
near — | Customers have rated it — | . |  
| ... | ... | ... | ... | ... | .

# Neural Template Generation Approach

## Step 1: Encode the Source

Fitzbillies,ty[coffee shop],pr[< £20],food[Chinese],cust[3/5],area[city centre]

## Step 2: Select a Template



## Step 3: Fill-in Each Segment

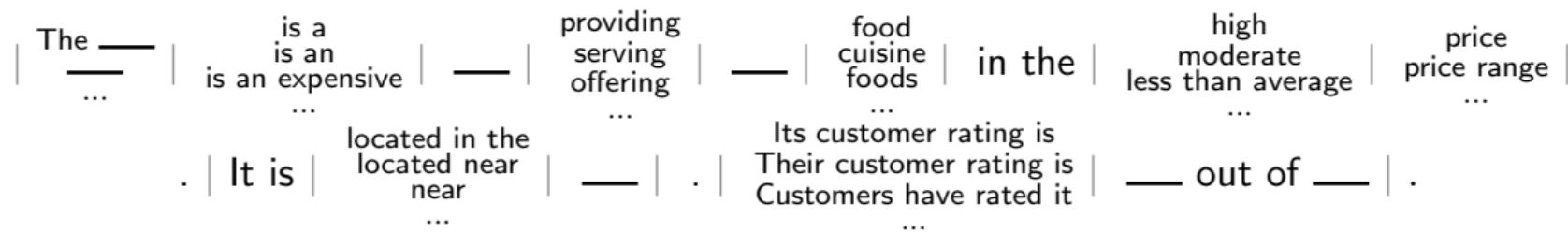
|| Fitzbillies ||

# Neural Template Generation Approach

## Step 1: Encode the Source

Fitzbillies,ty[coffee shop],pr[< £20],food[Chinese],cust[3/5],area[city centre]

## Step 2: Select a Template



## Step 3: Fill-in Each Segment

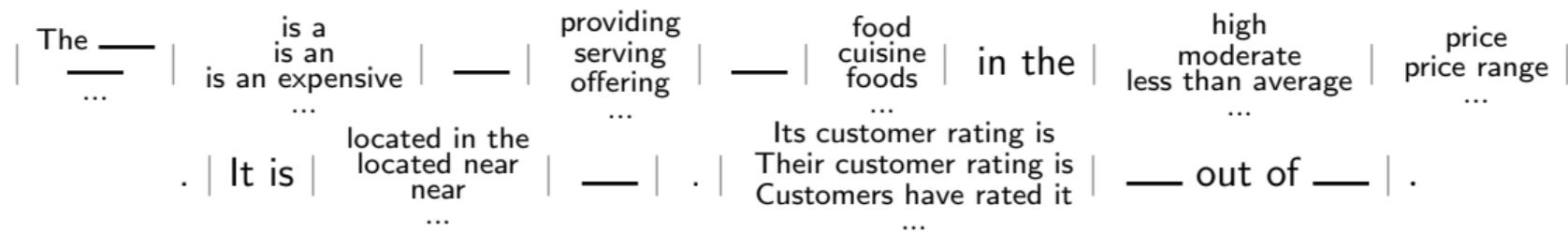
|| Fitzbillies || is a ||

# Neural Template Generation Approach

## Step 1: Encode the Source

Fitzbillies,ty[coffee shop],pr[< £20],food[Chinese],cust[3/5],area[city centre]

## Step 2: Select a Template



## Step 3: Fill-in Each Segment

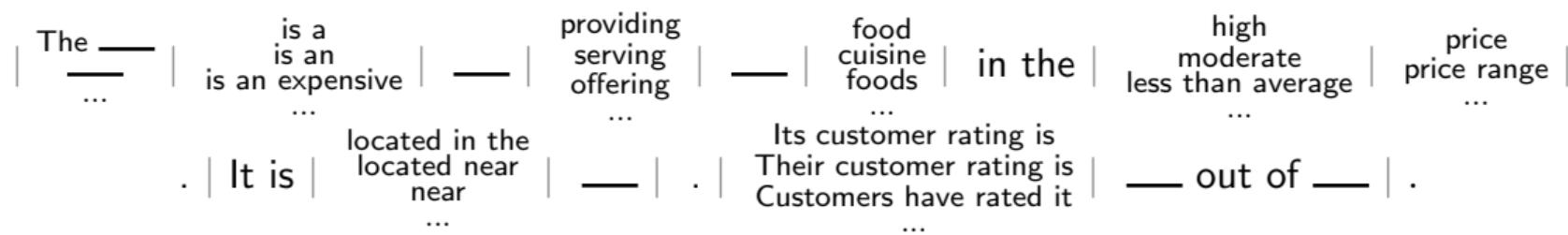
|| Fitzbillies || is a || coffee shop ||

# Neural Template Generation Approach

## Step 1: Encode the Source

Fitzbillies,ty[coffee shop],pr[< £20],food[Chinese],cust[3/5],area[city centre]

## Step 2: Select a Template

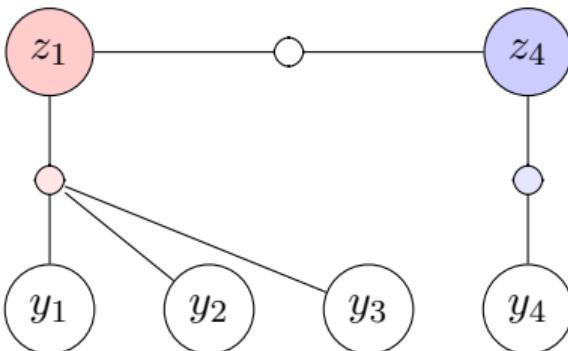


## Step 3: Fill-in Each Segment

|| Fitzbillies || is a || coffee shop || providing || Chinese || food || in the || moderate || price range || . || It is || located in the || city centre || . ||

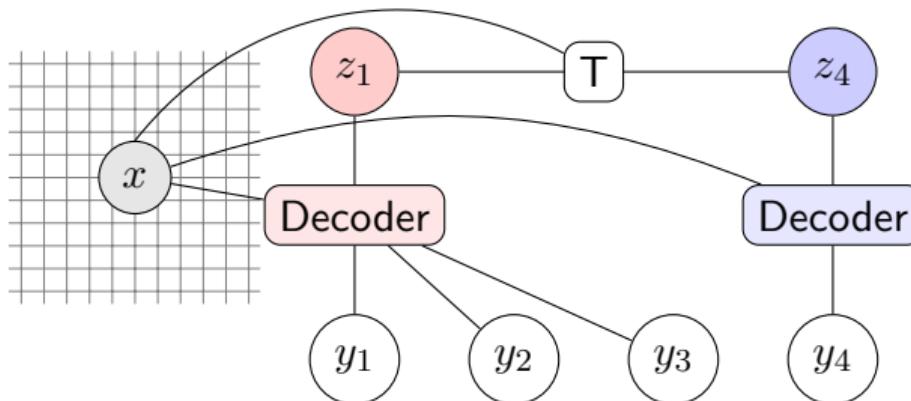
## Technical Methodology: Hidden Semi-Markov Model

- HMM: discrete latent states with single emissions (e.g. words).
- HSMM: discrete latent states produce multiple emissions (e.g. phrases).
- Parameterized with *transition*, *emission*, and *length* distributions.



# Technical Methodology: Neural Hidden Semi-Markov Model

- Employ HSMM as a conditional latent variable language model,  $p(y_1, \dots, y_T, z \mid x)$ .
- Transition Distribution: NN between states.
- Emission Distribution: Seq2Seq+Attention, one per state  $k$ .



## Technical Methodology: Learning Templates

- Fit model by maximizing log-marginal likelihood on training data.

$$\max_{\theta} \sum_j \log \sum_z p(y^{(j)}, z | x^{(j)}; \theta)$$

Details: Pre-score segmentations, HSMM forward algorithm for sum, backprop with autograd, all inference is exact.

## Technical Methodology: Learning Templates

- Fit model by maximizing log-marginal likelihood on training data.

$$\max_{\theta} \sum_j \log \sum_z p(y^{(j)}, z | x^{(j)}; \theta)$$

Details: Pre-score segmentations, HSMM forward algorithm for sum, backprop with autograd, all inference is exact.

- Compute argmax segmentations to find common *templates*.

$$z^{(j)} = \arg \max_z p(y^{(j)}, z | x^{(j)}; \theta)$$

[The Wrestlers]<sub>185</sub> [is a]<sub>29</sub> [coffee shop]<sub>164</sub> [that serves]<sub>188</sub> [English]<sub>139</sub> [food]<sub>18</sub> [in  
the]<sub>32</sub> [moderate]<sub>125</sub> [price range]<sub>180</sub> [.]<sub>90</sub>

# Neural Template

The — | is a — | providing — | food — | high — | price  
— | is an expensive | serving — | cuisine — | moderate — | price range  
... | ... | offering | foods | less than average | ...  
| ... | ... | ... | ... | ...  
| ... | ... | ... | ... | ... | ...  
. | It is | located in the — | Its customer rating is — | .  
| located near — | Their customer rating is — | .  
near | ... | ... | ... | ...  
| ... | ... | ... | ... | ... | .

# E2E Challenge

	BLEU	NIST
Test		
Substitution	43.78	6.88
Neural Template	56.72	7.63
Full Neural Model	65.93	8.59

	BLEU	NIST	ROUGE-4
Conditional KN-LM	19.8	5.19	10.7
NNLM (field)	33.4	7.52	23.9
NNLM (field & word)	34.7	7.98	25.8
Neural Template	33.8	7.51	28.2

# Issue 1: Interpretability

---

## kenny warren

---

**name:** kenny warren, **birth date:** 1 april 1946,

**birth name:** kenneth warren deutsch, **birth place:** brooklyn, new york,

**occupation:** ventriloquist, comedian, author,

**notable work:** book - the revival of ventriloquism in america

---

1. kenneth warren deutsch ( april 1, 1946 ) is an american ventriloquist.
  2. kenneth warren deutsch ( april 1, 1946 , brooklyn,) is an american ventriloquist.
  3. kenneth warren deutsch ( april 1, 1946 ) is an american  
ventriloquist, best known for his the revival of ventriloquism.
  4. “kenny” warren is an american ventriloquist.
  5. kenneth warren “kenny” warren (born april 1, 1946 ) is  
an american ventriloquist, and author.
-

## Issue 2: Controllability

---

### The Golden Palace

---

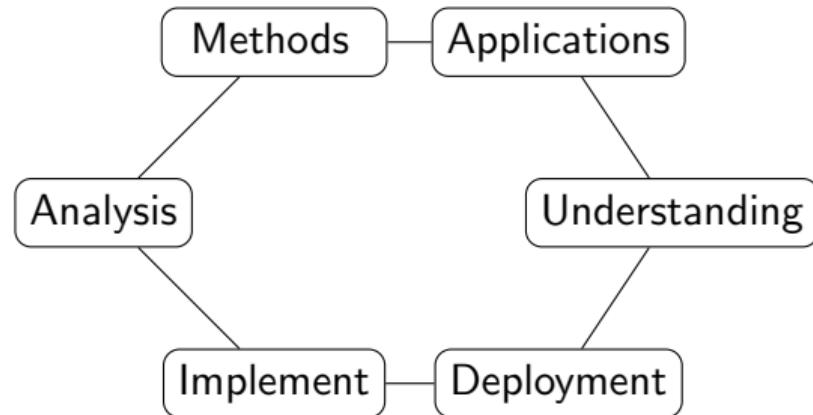
name[The Golden Palace], type[coffee shop], food[Chinese],  
priceRange[cheap] custRating[5 out of 5], area[city centre],

---

1. The Golden Palace is a coffee shop located in the city centre.
  2. In the city centre is a cheap Chinese coffee shop called  
The Golden Palace.
  3. The Golden Palace that serves Chinese food in the cheap  
price range. It is located in the city centre. Its customer  
rating is 5 out of 5.
  4. The Golden Palace is a Chinese coffee shop.
  5. The Golden Palace is a Chinese coffee shop  
with a customer rating of 5 out of 5.
-

# Future Work

NLP post deep learning



# Long-Form Generation with Explicit Reasoning

TEAM	WIN	LOSS	PTS	FG.PCT	RB	AS ...
Hawks	11	12	103	49	47	27
Heat	7	15	95	43	34	20

PLAYER	AS	RB	PT	FG	FGA	CITY ...
Tyler Johnson	5	2	27	8	16	Miami
Dwight Howard	11	17	23	9	11	Atlanta
Paul Millsap	2	9	21	8	12	Atlanta
Goran Dragic	4	2	21	8	17	Miami
Wayne Ellington	2	3	19	7	15	Miami
Dennis Schroder	7	4	17	8	15	Atlanta
Rodney McGruder	5	5	11	3	8	Miami
Hasan Whiteside	2	12	8	4	12	Miami

...

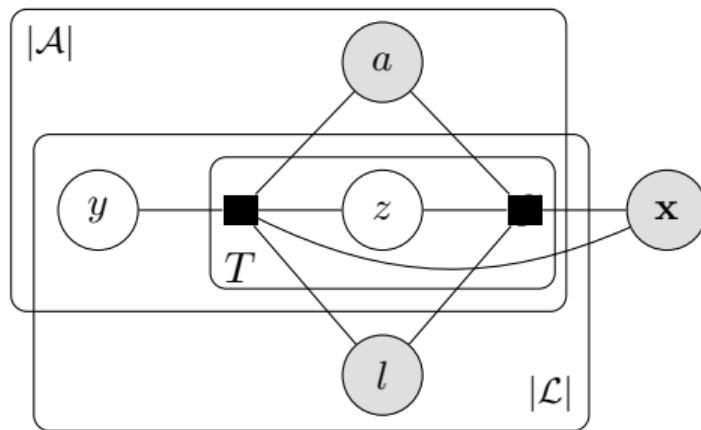
(2)

(1)

[The Atlanta Hawks defeated the Miami Heat, 103 - 95, at Philips Arena on Wednesday.] [Defense was key for the Hawks, as they held the Heat to 42 percent shooting and forced them to commit 16 turnovers. Atlanta also dominated in the paint, winning the rebounding battle, 47 - 34, and outscoring them in the paint 58 - 26. The Hawks shot 49 percent from the field and assisted on 27 of their 43 made baskets.] [Miami ( 7 - 15 ) are as beat-up as anyone right now. Hassan Whiteside really struggled in this game, as he amassed eight points, 12 rebounds and one blocks on 4 of - 12 shooting] ...

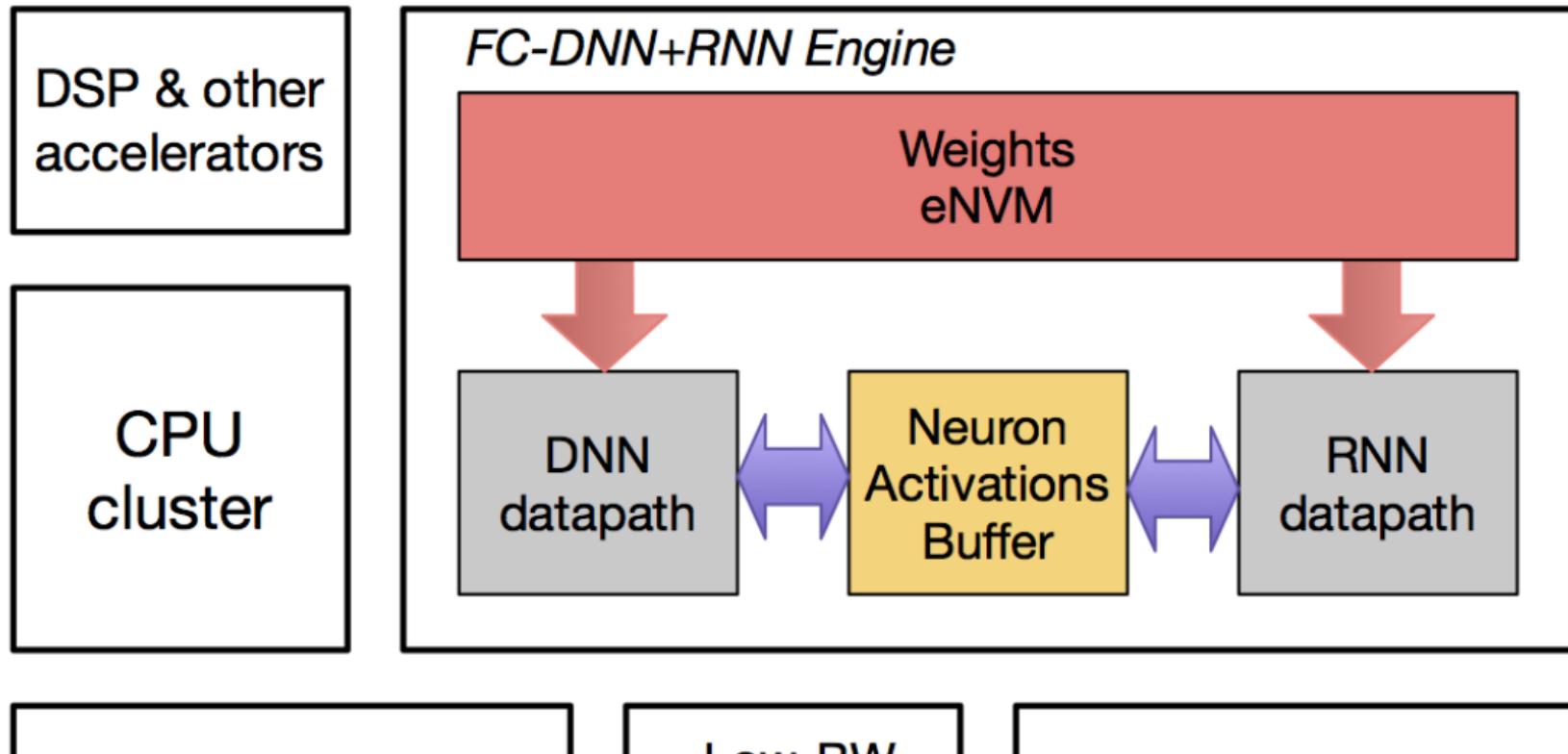
# Probabilistic Programming

(Preprint)



# Learning Neural Reasoning-Based Models

## *Universal Translator SoC*





Yuntian Deng, Anssi Kanervisto, and Alexander M. Rush. 2016. What You Get Is What You See: A Visual Markup Decompiler. In *Arxiv*.

Yuntian Deng, Yoon Kim, Justin Chiu, Demi Guo, and Alexander Rush. 2018. Latent alignment and variational attention. In *Advances in Neural Information Processing Systems*, pages 9735–9747.

Yoon Kim, Carl Denton, Luong Hoang, and Alexander M. Rush. 2017. Structured attention networks. [abs/1702.00887](#).

Yoon Kim, Yacine Jernite, David Sontag, and Alexander M. Rush. 2016. Character-Aware Neural Language Models. In *AAAI*.

Yoon Kim and Alexander M. Rush. 2016. Sequence-Level Knowledge Distillation. In *EMNLP*.

Yoon Kim, Sam Wiseman, Andrew C. Miller, David Sontag, and Alexander M. Rush. 2018. Semi-amortized variational autoencoders.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush. 2017.

Opennmt: Open-source toolkit for neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, System Demonstrations*, pages 67–72.

Brandon Reagen, Udit Gupta, Robert Adolf, Michael M Mitzenmacher, Alexander M Rush, Gu-Yeon Wei, and David Brooks. 2017. Weightless: Lossy weight encoding for deep neural network compression. *arXiv preprint arXiv:1711.04686*.

Alexander Rush. 2018. The annotated transformer. In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 52–60.

Alexander M Rush, Sumit Chopra, and Jason Weston. 2015. A Neural Attention Model for Abstractive Sentence Summarization. In *EMNLP*, September, pages 379–389.

Allen Schmaltz, Yoon Kim, Alexander M. Rush, and Stuart M. Shieber. 2016. Sentence-Level Grammatical Error Identification as Sequence-to-Sequence Correction. In *arxiv*.

Jean Senellart, Dakun Zhang, WANG Bo, Guillaume Klein, Jean-Pierre Ramatchandrin,

- Josep Crego, and Alexander Rush. 2018. Opennmt system description for wnmnt 2018: 800 words/sec on a single-core cpu. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 122–128.
- Hendrik Strobelt, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, and Alexander M Rush. 2019. Seq2seq-v is: A visual debugging tool for sequence-to-sequence models. *IEEE transactions on visualization and computer graphics*, 25(1):353–363.
- Hendrik Strobelt, Sebastian Gehrmann, Bernd Huber, Hanspeter Pfister, and Alexander M. Rush. 2016. Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks. In *Arxiv*.
- Sam Wiseman and Alexander M. Rush. 2016. Sequence-to-Sequence Learning as Beam-Search Optimization. In *EMNLP*.
- Sam Wiseman, Alexander M. Rush, and Stuart M. Shieber. 2017a. Challenges in Data-to-Document Generation. In *EMNLP*.

Sam Wiseman, Stuart M. Shieber, and Alexander M. Rush. 2017b. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 2253–2263.

Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2018. Learning neural templates for text generation. *arXiv preprint arXiv:1808.10122*.