# The Amazing World of Neural Language Generation

## Part II: Neural Network Modeling for Generation

---

Yangfeng Ji

November 20, 2020

Department of Computer Science
University of Virginia

UNIVERSITY *of* VIRGINIA | ENGINEERING

# Basic Architecture of Neural NLG Models

From a (very) high-level viewpoint, neural NLG model can be formulated as with two fundamental components: **encoder** and **decoder**:

Input $x$ ⟶ Encoder ⟶ Decoder ⟶ Text $y$

where

- ▶ Input: A sequence of words $x = (x_1, \cdots, x_m)$, $m$ words
- ▶ Output: A sequence of words $y = (y_1, \cdots, y_n)$, $n$ words

# Basic Architecture of Neural NLG Models

From a (very) high-level viewpoint, neural NLG model can be formulated as with two fundamental components: **encoder** and **decoder**:
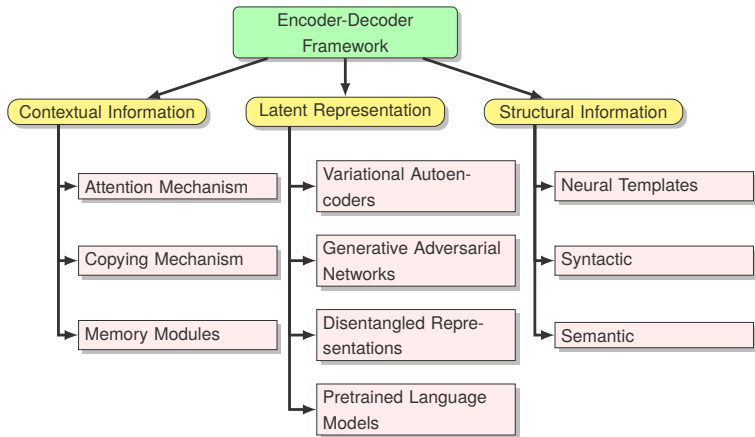
$$\text{Input } x \longrightarrow \boxed{\text{Encoder}} \longrightarrow \boxed{\text{Decoder}} \longrightarrow \text{Text } y$$

where

- Input: A sequence of words $x = (x_1, \cdots, x_m)$, $m$ words
- Output: A sequence of words $y = (y_1, \cdots, y_n)$, $n$ words

We will try to answer the following three questions with neural network modeling strategies
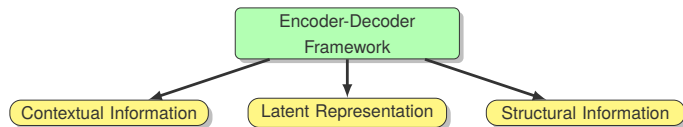
- How to select *contextual information*?
- How to build better *latent representations*?
- How to incorporate *structural information*?

This part of the tutorial will cover the three major neural network modeling strategies for text generation

# Roadmap
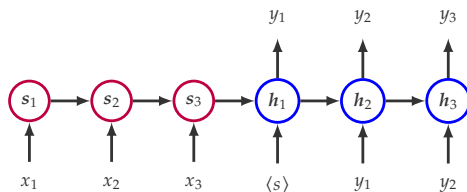
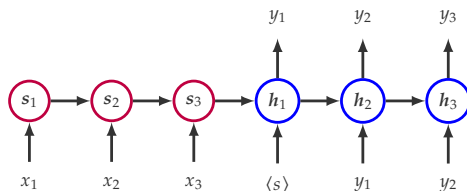This part of the tutorial will cover the three major neural network modeling strategies for text generation

The simple implementation of the encoder-decoder framework is to realize each component with a recurrent neural network as illustrated in the following



where each $s_.$/$h_.$ is a hidden state of the encoder/decoder recurrent neural network.

# Recurrent Neural Networks as Encoders/Decoders

The simple implementation of the encoder-decoder framework is to realize each component with a recurrent neural network as illustrated in the following



where each $s_.$/$h_.$ is a hidden state of the encoder/decoder recurrent neural network.
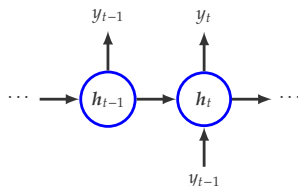
Simple extensions on the encoder side

- ▶ Bi-directional RNN
- ▶ Stacked LSTM

# Recurrent Neural Network Decoder

In general, an decoder can be implemented as an auto-regressive model, with the hidden state computed as
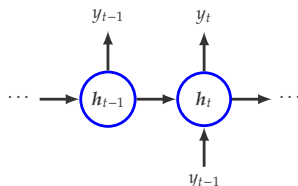
$$h_t = f(h_{t-1}, y_{t-1}) \tag{1}$$

# Recurrent Neural Network Decoder

In general, an decoder can be implemented as an auto-regressive model, with the hidden state computed as

$$h_t = f(h_{t-1}, y_{t-1}) \tag{1}$$
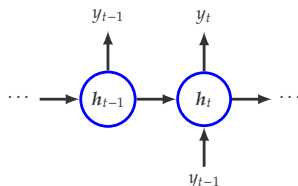


For generation, the probability of the word $y_t$ is computed as

$$p(y) = \text{softmax}(W_o h_t) \tag{2}$$

where $W_o \in \mathbb{R}^{H \times V}$ is a learnable weight matrix for the output layer

# Recurrent Neural Network Decoder

In general, an decoder can be implemented as an auto-regressive model, with the hidden state computed as

$$h_t = f(h_{t-1}, y_{t-1}) \tag{1}$$



For generation, the probability of the word $y_t$ is computed as

$$p(y) = \text{softmax}(W_o h_t) \tag{2}$$
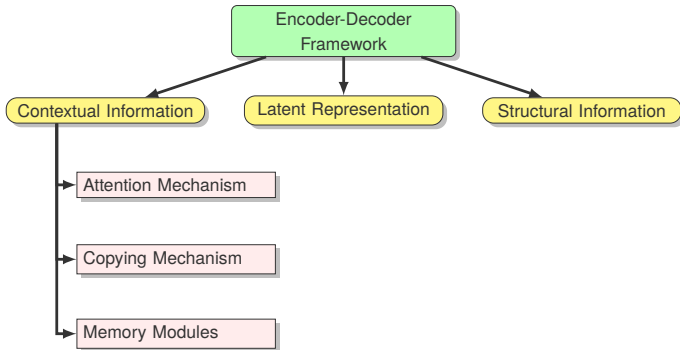
where $W_o \in \mathbb{R}^{H \times V}$ is a learnable weight matrix for the output layer

Decoding algorithms will be discussed in the next part of the tutorial.
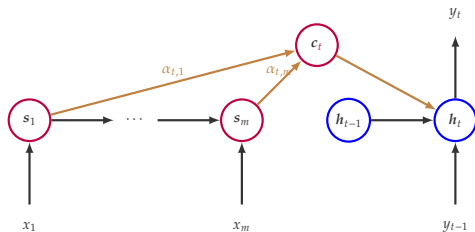
# Contextual Information

Attention mechanism (Bahdanau et al., 2015) provides a way of actively encoding context information from preceding text $x = (x_1, \cdots, x_m)$ via hidden states $(s_1, \ldots, s_m)$
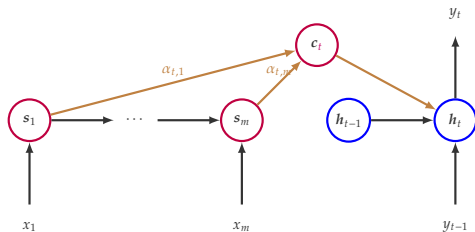
# Attention Mechanism

Attention mechanism (Bahdanau et al., 2015) provides a way of actively encoding context information from preceding text $x = (x_1, \cdots, x_m)$ via hidden states $(s_1, \ldots, s_m)$



Attention weights

$$\alpha_{t,i} = g(s_i, h_{t-1}) \in (0, 1)$$

$$c_t = \sum_{i=1}^{m} \alpha_{t,i} \cdot s_i$$

# Attention Mechanism

Attention mechanism (Bahdanau et al., 2015) provides a way of actively encoding context information from preceding text $x = (x_1, \cdots, x_m)$ via hidden states $(s_1, \ldots, s_m)$
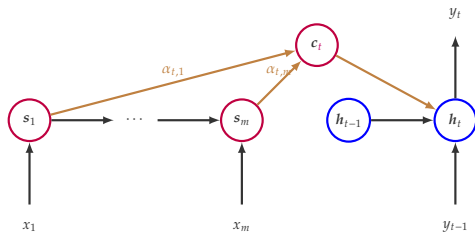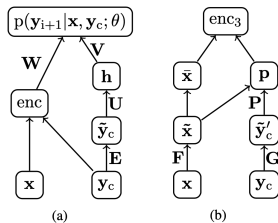


Attention weights

$$\alpha_{t,i} = g(s_i, h_{t-1}) \in (0,1)$$

$$c_t = \sum_{i=1}^{m} \alpha_{t,i} \cdot s_i$$

- ▶ $c_t$ is a function of $h_{t-1}$, which means it dynamically changes at every time step of decoding
- ▶ $c_t$ enables the decoder to be more selective on using contextual information
- ▶ Widely used in many text generation tasks: response generation, story generation, and document summarization

Rush et al. (2015) use the attention mechanism in a feed-forward neural network for abstractive sentence summarization, where the neural network architecture is constructed as the following



(a)                    (b)

Rush et al. (2015) use the attention mechanism in a feed-forward neural network for abstractive sentence summarization, where the neural network architecture is constructed as the following
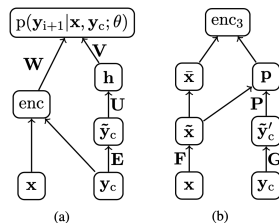


(a)        (b)

Diagram (b) represents the attention-based encoder, which use the input $x$ and a fixed context window $y_c = y_{i-c+1:i}$ to compute the attention weights
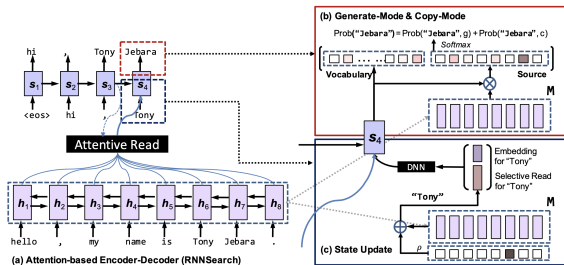
$$\boldsymbol{\alpha} \propto \exp(\tilde{x}P y'_c) \quad \bar{x}_i = \sum_{q=i-Q}^{i+Q} \tilde{x}_i/Q \quad \text{enc}(x, y_c) = \boldsymbol{\alpha}^\top \bar{x} \qquad (3)$$

Gu et al. (2016) propose a model called COPYNET to directly copy a phrases from input $x$ to output $y$,

$$p(y_t \mid h_t) = p_g(y_t \mid h_t) + p_c(y_t \mid h_t) \qquad (4)$$

where $p_g(y_t \mid h_t)$ is a probability distribution defined on $\mathcal{V}$ and $p_c(y_t \mid h_t)$ is a probability distribution defined only on the input $x$.
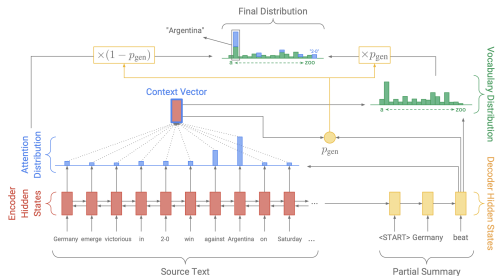


(a) Attention-based Encoder-Decoder (RNNSearch)

In other words, it defines a mixture model with equal mixture coefficients.

# Copying Mechanism: Pointer-Generator Networks

See et al. (2017) propose a similar idea to copy words from input $x$ to output $y$ in text summarization. The probability of $w \in \mathcal{Y}$ being the next word is

$$p(y_t = w \mid \cdot) = \beta p_g(y_t = w \mid \cdot) + (1 - \beta) \sum_{x_i \in x} \delta(w = x_i) \alpha_{t,i} \qquad (5)$$

Among the many implementation differences with (Gu et al., 2016), the work uses

# Copying Mechanism: Pointer-Generator Networks

See et al. (2017) propose a similar idea to copy words from input $x$ to output $y$ in text summarization. The probability of $w \in \mathcal{Y}$ being the next word is

$$p(y_t = w \mid \cdot) = \beta p_g(y_t = w \mid \cdot) + (1 - \beta) \sum_{x_i \in x} \delta(w = x_i)\alpha_{t,i} \qquad (5)$$

Among the many implementation differences with (Gu et al., 2016), the work uses

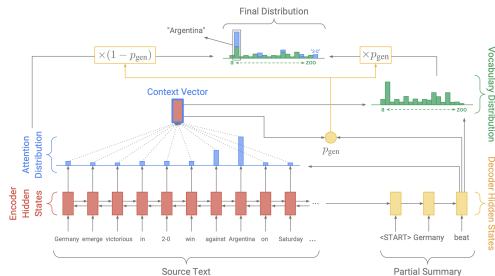▶ a soft gate $\beta \in (0, 1)$ to decide the probability of generation instead of copying
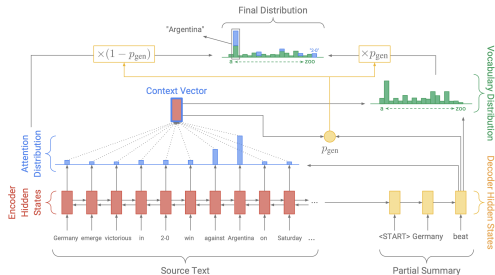
# Copying Mechanism: Pointer-Generator Networks

See et al. (2017) propose a similar idea to copy words from input $x$ to output $y$ in text summarization. The probability of $w \in \mathcal{Y}$ being the next word is

$$p(y_t = w \mid \cdot) = \beta p_g(y_t = w \mid \cdot) + (1 - \beta) \sum_{x_i \in x} \delta(w = x_i) \alpha_{t,i} \qquad (5)$$
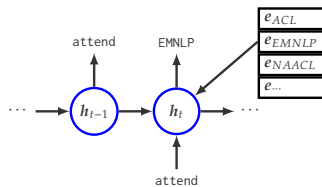
Among the many implementation differences with (Gu et al., 2016), the work uses

► a soft gate $\beta \in (0, 1)$ to decide the probability of generation instead of copying

► the probability of copy a word $w$ is defined by the attention weights associated with it

# Memory Modules

One example is proposed in (Clark et al., 2018) for entity-driven text (story) generation, where each memory cell is associated with a particular entity, to encode entity related information from context



where $e_\cdot$ is a distributed representation of an entity
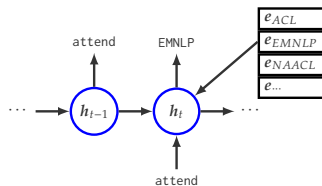
# Memory Modules

One example is proposed in (Clark et al., 2018) for entity-driven text (story) generation, where each memory cell is associated with a particular entity, to encode entity related information from context



where $e.$ is a distributed representation of an entity

- Entity prediction $p(e = \text{EMNLP}) \propto \exp(h_{t-1}^\top W_e e_{\text{EMNLP}} + w_f f(e))$, where $f$ is the surface feature related to this entity (Ji et al., 2017).

# Memory Modules

One example is proposed in (Clark et al., 2018) for entity-driven text (story) generation, where each memory cell is associated with a particular entity, to encode entity related information from context
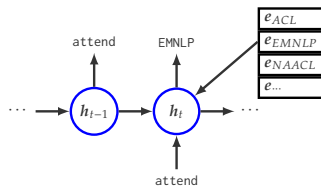


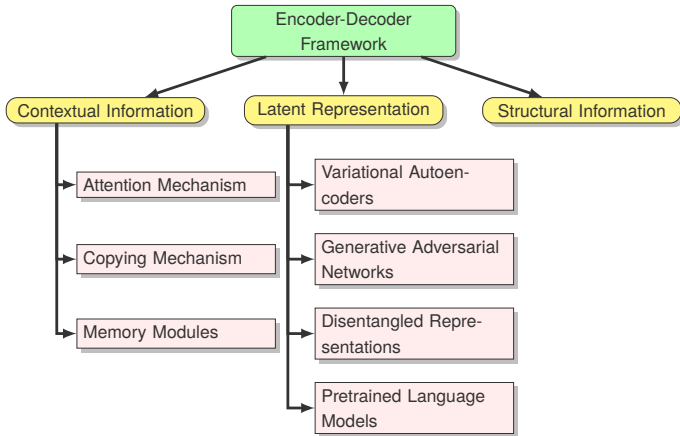where $e_.$ is a distributed representation of an entity

- Entity prediction $p(e = \text{EMNLP}) \propto \exp(h_{t-1}^\top W_e e_{\text{EMNLP}} + w_f f(e))$, where $f$ is the surface feature related to this entity (Ji et al., 2017).
- Dynamic entity updating $e_{\text{EMNLP}}^{(\text{new})} \propto \delta_t e_{\text{EMNLP}}^{(\text{new})} + (1 - \delta_t) h_t$, where $\delta_t \in (0, 1)$ determines how much information should be encoded from $e_{\text{EMNLP}}^{(\text{new})}$.

# Latent Representation

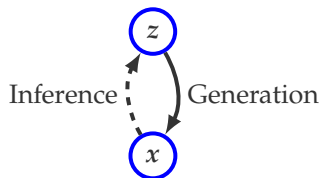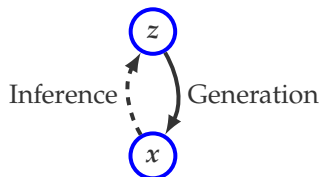One-page summary of variational autoencoder (Kingma and Welling, 2014)

# Variational Autoencoders in NLG

One-page summary of variational autoencoder (Kingma and Welling, 2014)



A formulation of variational autoencoders for text generation

$$h = \text{Encoder}(x) \qquad\qquad (6)$$

$$z = h + \epsilon \quad \epsilon \sim \mathcal{N}(0, \text{diag}(\sigma^2)) \qquad (7)$$

$$\tilde{x} = \text{Decoder}(z) \qquad\qquad (8)$$

# Variational Autoencoders in NLG

One-page summary of variational autoencoder (Kingma and Welling, 2014)



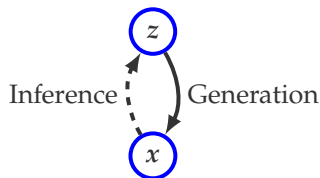A formulation of variational autoencoders for text generation

$$h = \text{Encoder}(x) \tag{6}$$
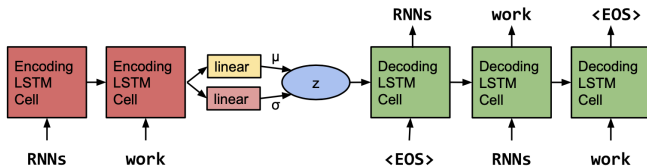
$$z = h + \epsilon \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \text{diag}(\sigma^2)) \tag{7}$$

$$\tilde{x} = \text{Decoder}(z) \tag{8}$$

An impact of using VAE is that it (1) produces a robust encoder for input $x$ and (2) enriches the hidden space $\mathcal{H}$.

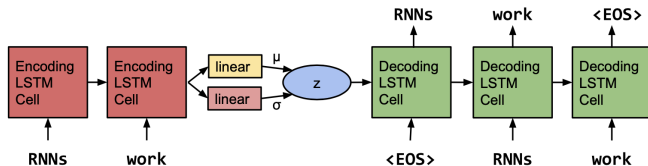# Variational Seq2seq Models

An example application of variational autoencoder in language generation



(Bowman et al., 2016)

# Variational Seq2seq Models

An example application of variational autoencoder in language generation



- ▶ The mean and variance of latent variable $z$ is computed by the linear transformations of the last hidden states from the RNN encoder $s_m$

(Bowman et al., 2016)

# Variational Seq2seq Models

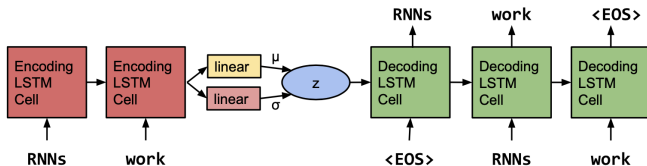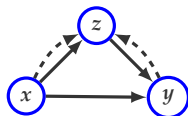An example application of variational autoencoder in language generation



▶ The mean and variance of latent variable $z$ is computed by the linear transformations of the last hidden states from the RNN encoder $s_m$

▶ Other influential ideas from this work are *KL cost annealing* and *adversarial evaluation*

(Bowman et al., 2016)

# Conditional VAE

A simple formulation of conditional VAE is proposed by Sohn et al. (2015), which initially was used in computer vision.



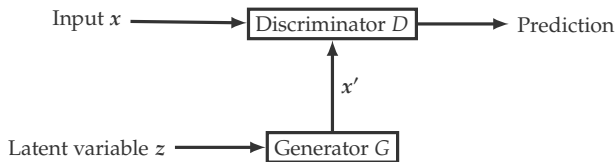In text generation, consider $x$, $y$ and $z$ are input texts, output texts, and latent representations of input texts

- ▶ Generation network: $p_\theta(y \mid x, z)$, where $z \sim p_\theta(z \mid x)$
- ▶ Inference network: $q_\phi(z \mid x, y)$
- ▶ *Example applications*: in response generation, latent variable $z$ indicates whether the next utterance is a generic response (Shen et al., 2017)
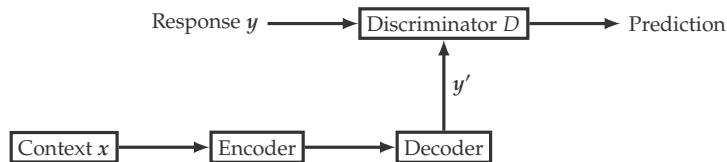
The basic pipeline of GAN is described in the following pipeline



- ▶ One goal is to learn a generator $G$ that can generate text $x'$ with the same quality as $x$ (in other words, to "fool" the discriminator)
- ▶ Potential applications are to adopt the framework as one component in other task-specific generation tasks (e.g., style transfer)

# GANs for Generation

As a straightforward application of adversarial learning is to replace the generator with a sequence-to-sequence model as we discussed before, as proposed by Li et al. (2017)



- ▶ The decoder is to generate a response $y'$ based on input context $x$
- ▶ The discriminator is to predict whether a response is generated by humans or the decoder

# Learning Disentangled Representations

One way to utilize a discriminator is to learn disentangled representations and make sure latent representations encoding the *expected attributes* for generation.

# Learning Disentangled Representations

One way to utilize a discriminator is to learn disentangled representations and make sure latent representations encoding the *expected attributes* for generation.

An example of learning disentangled representations is proposed by Hu et al. (2017), with

▶ Encoder
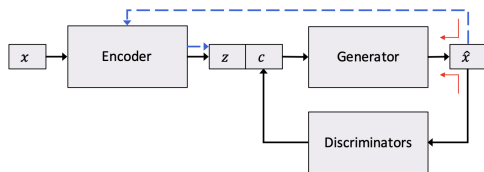  $(z) = \text{Encoder}(x)$

▶ Decoder
  $\hat{x} = \text{Decoder}(z, c)$

▶ Discriminator
  $c = \text{Dis}(\hat{x})$



where $c$ encodes the attributes of a text (e.g., sentiment categories, formality). Other work on learning disentangled representations include (Bao et al., 2019; Cheng et al., 2020)

Vaswani et al. (2017): "*Attention mechanism . . . , allowing modeling of dependencies without regard to their distance in the input or output sequences.*"



Recent applications of transformers simply use them as basic building blocks, just like the way of using LSTM

# Generative Pre-trained Transformers (GPT)

GPT (Radford et al., 2018) is trained simply by predicting the next words with

$$h_0 = W_e x_{t-k:t-1} + W_p$$
$$h_l = \text{transformer\_block}(h_{l-1}) \; \forall l \in [1, 12] \tag{9}$$
$$p(y_t \mid h_n) = \text{softmax}(W_o h_n)$$

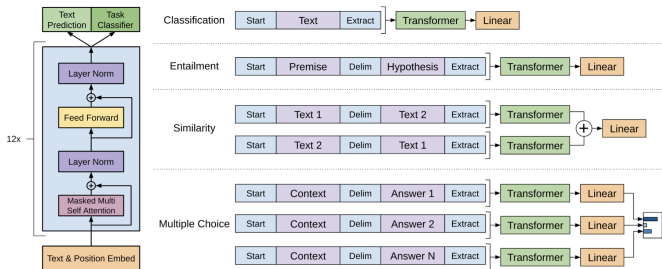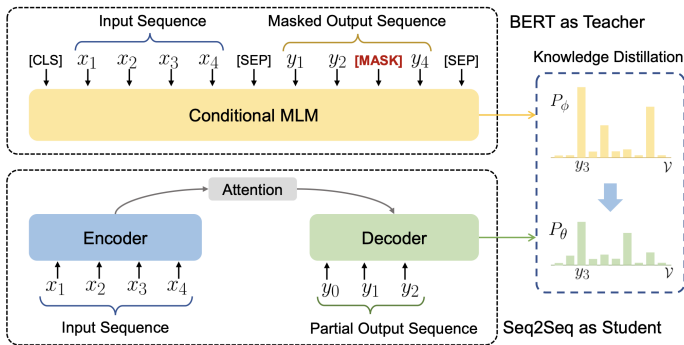where $W_e$ and $W_p$ are the word and position embeddings.



Figure 1: **(left)** Transformer architecture and training objectives used in this work. **(right)** Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

# Pre-trained Models as Teachers

Besides the generative models like GPT-2 and its variants, we can also use BERT for text generation. One example from (Chen et al., 2020b) is that BERT can be used as a teacher model to help train a sequence-to-sequence models for better output probability
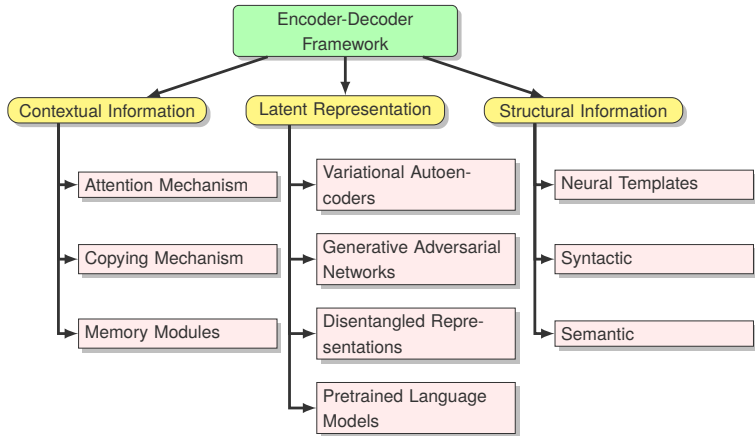
# Structural Information

# Generation with Neural Templates

An extension of variational antoencoders is to incorporate sequential information in latent variables. For example, Wiseman et al. (2018) propose a semi-Markov model on the latent variable sequence $z = (z_1, \ldots, z_n)$ to capture dependency among adjacent words for text-to-data generation.

**Source Entity**: Cotto
type[coffee shop], rating[3 out of 5], food[English], area[city centre], price[moderate], near[The Portland Arms]

**System Generation:**
Cotto is a coffee shop serving English food in the moderate price range. It is located near The Portland Arms. Its customer rating is 3 out of 5.

**Neural Template:**



A neural hidden semi-Markov model decoder

▶ Transition distribution $p(z_{t+1} \mid z_t, x)$

▶ Length distribution $p(l_{t+1} \mid z_{t+1})$

▶ Emission distribution
$p(y_{t-l_t+1:t} \mid z_t = k, l_t = l, x)$

# Graph-to-Sequence Generation

An example of inputs is an AMR graph, where the task is to generate a text with the same meaning as the input graph.

An example of inputs is an AMR graph, where the task is to generate a text with the same meaning as the input graph.



▶ The whole model is a layer-wise LSTM, within each layer LSTM runs on the linearized graph, and between layers the hidden states are updated with neighbor nodes from the previous layer

# Graph-to-Sequence Generation

An example of inputs is an AMR graph, where the task is to generate a text with the same meaning as the input graph.



- ▶ The whole model is a layer-wise LSTM, within each layer LSTM runs on the linearized graph, and between layers the hidden states are updated with neighbor nodes from the previous layer
- ▶ The idea of linearizing an AMR graph (using depth-first search) is proposed by Konstas et al. (2017)
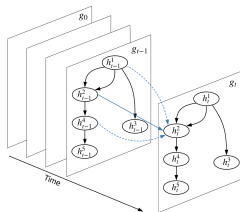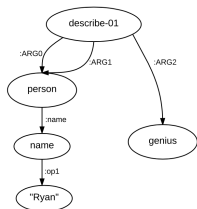
# Graph-to-Sequence Generation

An example of inputs is an AMR graph, where the task is to generate a text with the same meaning as the input graph.



- ▶ The whole model is a layer-wise LSTM, within each layer LSTM runs on the linearized graph, and between layers the hidden states are updated with neighbor nodes from the previous layer
- ▶ The idea of linearizing an AMR graph (using depth-first search) is proposed by Konstas et al. (2017)
- ▶ Other ideas of incorporating AMR graphs: Graph neural network (Chen et al., 2020a)

# Dialogue Generation with Semantic Exemplars

A straightforward application of GPT is the response generation proposed by Gupta et al. (2020). The prediction is still done by word-by-word prediction, where the input sequence and output sequence are [1]

$$x = (\text{Dialogue context}, \text{Response frames}, \text{Response text})$$
$$y = (\langle \text{MASK} \rangle, \text{Response frames}, \text{Response text})$$

$$(10)$$



---

[1]A semantic frame example PERCEPTION: hear, say, see, smell, feel

# Reference I

Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural Machine Translation by Jointly Learning to Align and Translate. In *International Conference on Learning Representations*. arXiv: 1409.0473.

Bao, Y., Zhou, H., Huang, S., Li, L., Mou, L., Vechtomova, O., Dai, X.-y., and Chen, J. (2019). Generating Sentences from Disentangled Syntactic and Semantic Spaces. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6008–6019, Florence, Italy. Association for Computational Linguistics.

Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A., Jozefowicz, R., and Bengio, S. (2016). Generating Sentences from a Continuous Space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21, Berlin, Germany. Association for Computational Linguistics.
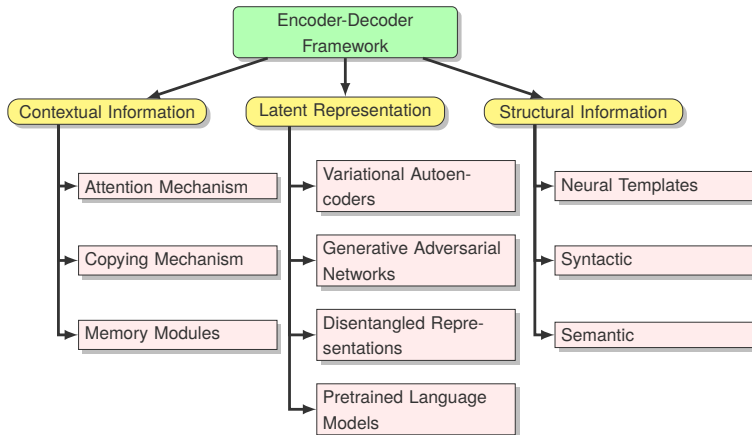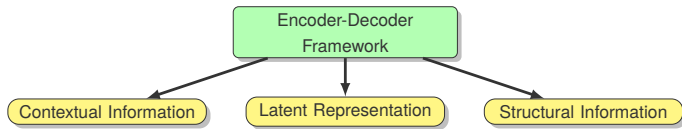
Chen, M., Tang, Q., Wiseman, S., and Gimpel, K. (2019). Controllable Paraphrase Generation with a Syntactic Exemplar. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5972–5984, Florence, Italy. Association for Computational Linguistics.

Chen, Y., Wu, L., and Zaki, M. J. (2020a). Reinforcement Learning Based Graph-to-Sequence Model for Natural Question Generation. *arXiv:1908.04942 [cs]*. arXiv: 1908.04942.

Chen, Y.-C., Gan, Z., Cheng, Y., Liu, J., and Liu, J. (2020b). Distilling Knowledge Learned in BERT for Text Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7893–7905, Online. Association for Computational Linguistics.

Cheng, P., Min, M. R., Shen, D., Malon, C., Zhang, Y., Li, Y., and Carin, L. (2020). Improving Disentangled Text Representation Learning with Information-Theoretic Guidance. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7530–7541, Online. Association for Computational Linguistics.

Clark, E., Ji, Y., and Smith, N. A. (2018). Neural Text Generation in Stories Using Entity Representations as Context. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2250–2260, New Orleans, Louisiana. Association for Computational Linguistics.

Gu, J., Lu, Z., Li, H., and Li, V. O. (2016). Incorporating Copying Mechanism in Sequence-to-Sequence Learning. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1631–1640, Berlin, Germany. Association for Computational Linguistics.

Gupta, P., Bigham, J. P., Tsvetkov, Y., and Pavel, A. (2020). Controlling Dialogue Generation with Semantic Exemplars. *arXiv:2008.09075 [cs]*. arXiv: 2008.09075.

Hu, Z., Yang, Z., Liang, X., Salakhutdinov, R., and Xing, E. P. (2017). Toward Controlled Generation of Text. volume 70 of *Proceedings of Machine Learning Research*, pages 1587–1596, International Convention Centre, Sydney, Australia. PMLR.

# Reference II

Ji, Y., Tan, C., Martschat, S., Choi, Y., and Smith, N. A. (2017). Dynamic Entity Representations in Neural Language Models. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1830–1839, Copenhagen, Denmark. Association for Computational Linguistics.

Kingma, D. P. and Welling, M. (2014). Auto-Encoding Variational Bayes. In *ICLR*. arXiv: 1312.6114.

Konstas, I., Iyer, S., Yatskar, M., Choi, Y., and Zettlemoyer, L. (2017). Neural AMR: Sequence-to-Sequence Models for Parsing and Generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 146–157, Vancouver, Canada. Association for Computational Linguistics.

Li, J., Monroe, W., Shi, T., Jean, S., Ritter, A., and Jurafsky, D. (2017). Adversarial Learning for Neural Dialogue Generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2157–2169, Copenhagen, Denmark. Association for Computational Linguistics.

Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). *Improving language understanding by generative pre-training*.

Rush, A. M., Chopra, S., and Weston, J. (2015). A Neural Attention Model for Abstractive Sentence Summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 379–389, Lisbon, Portugal. Association for Computational Linguistics.

See, A., Liu, P. J., and Manning, C. D. (2017). Get To The Point: Summarization with Pointer-Generator Networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.

Shen, X., Su, H., Li, Y., Li, W., Niu, S., Zhao, Y., Aizawa, A., and Long, G. (2017). A Conditional Variational Framework for Dialog Generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 504–509, Vancouver, Canada. Association for Computational Linguistics.

Sohn, K., Lee, H., and Yan, X. (2015). Learning Structured Output Representation using Deep Conditional Generative Models. In *Advances in neural information processing systems*, pages 3483–3491.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is All you Need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.

Wiseman, S., Shieber, S., and Rush, A. (2018). Learning Neural Templates for Text Generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3174–3187, Brussels, Belgium. Association for Computational Linguistics.