

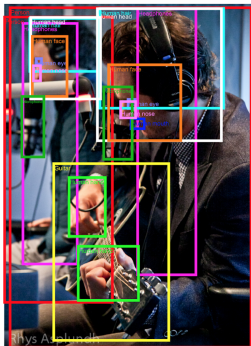
11-695: AI Engineering

Self-Supervision Introduction

LTI/SCS

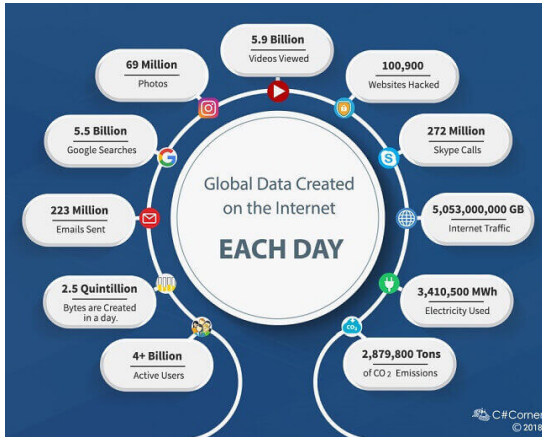
Spring 2020

- ① Data Issue
- ② Methodology
- ③ Auto Encoder
- ④ Data Encoding to Embedded Representation
- ⑤ Decoding from Embedded Space



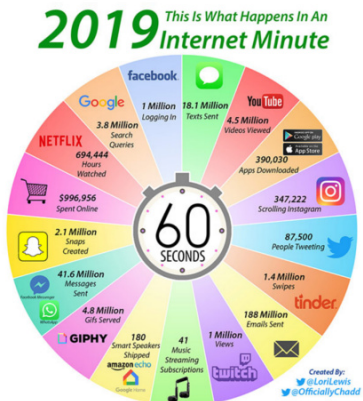
- Data is often carefully chosen,
- carefully cleaned and labelled,
- and usually *balanced*

But, The World's Data is Enormous



- Data is fast increased every second

Most Data is Unlabelled



- Also unstructured and unclean
- The cost to label them is prohibitively expensive.

Image credit: @LoriLewis & @OfficiallyChadd - Twitter

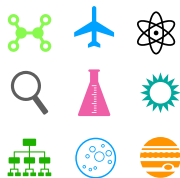
- ① Data Issue
- ② Methodology
- ③ Auto Encoder
- ④ Data Encoding to Embedded Representation
- ⑤ Decoding from Embedded Space



Raw Data

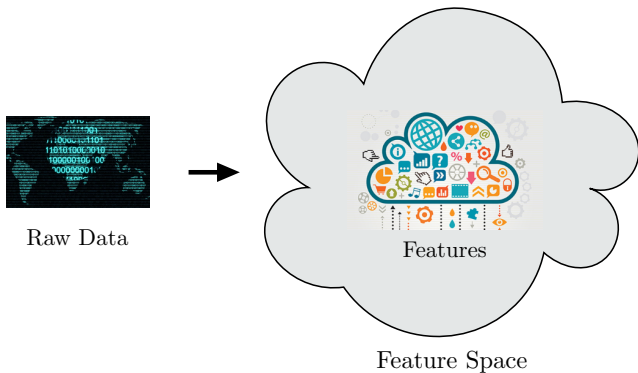


Model-based
Features

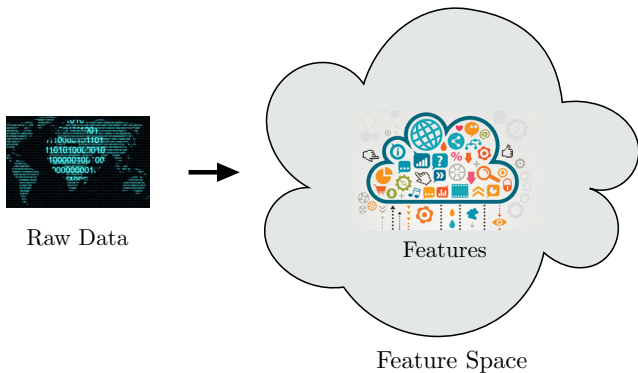


Tasks

- How about we not only focus on the task, but also put a proper focus on Data itself
- Difference from:
 - Fully labelled data: Supervised
 - Partially labelled data: Semi-supervised
- We have no labels thus no supervision → Self-supervised (Unsupervised).



- Classification/Regression: use features to predict labels/numbers
- Required: ground-truth labels/numbers

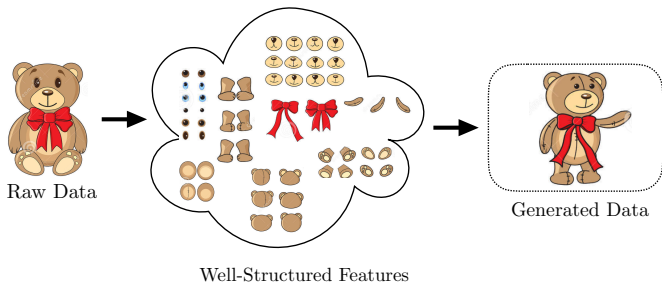


- Classification/Regression: use features to predict labels/numbers
- Required: ground-truth labels/numbers
- Unfortunately, we don't have such ground truths



- Use each data sample itself as a label
 - Extract Features from x : $z = f_1(x)$
 - Reconstruct data \hat{x} : $\hat{x} = f_2(z)$
 - Minimize the distance $d(x, \hat{x})$, *a.k.a* reconstruction error.
- But, data itself are often highly-dimensional \rightarrow (very) **hard!**
- Inject prior human knowledge/expertise about data, *e.g.*
 - Vision: use CNN as feature extractor
 - Text, sequential data: use RNN instead

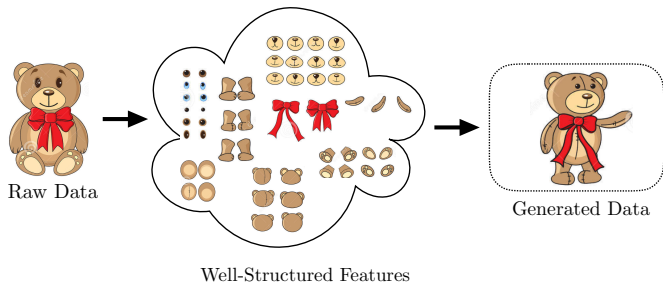
Generate a bear to just classify a bear?



- No
 - No in cases such as classification
 - In many cases, yes, we need to actually generate data
- Regardless the task, we need to have good data features¹
- Remind: for supervised, we don't need to generate data
 - Features needn't be well-structured to predict labels

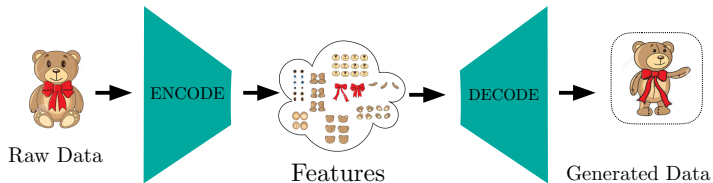
¹<https://arxiv.org/pdf/1206.5538.pdf>

Generate a bear to just classify a bear?

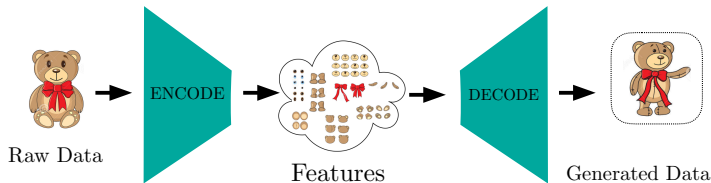


- Remind: for supervised, we don't need to generate data
 - Features needn't be well-structured to predict labels
- Comparison to Unsupervised in the classification/regression:
 - Supervised: we model $p(y|x)$
 - Unsupervised: $p(x, y)$

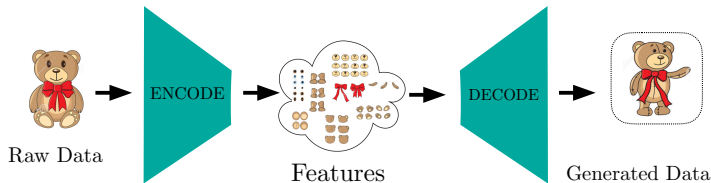
- ① Data Issue
- ② Methodology
- ③ Auto Encoder**
- ④ Data Encoding to Embedded Representation
- ⑤ Decoding from Embedded Space



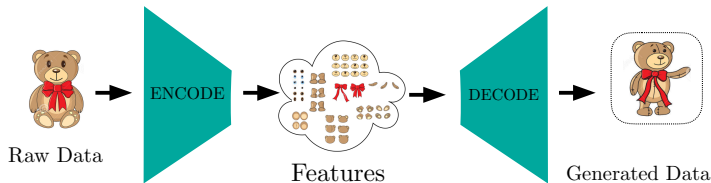
- Extract and Reconstruct *itself* \rightarrow *Auto-Encoder* + Decoder
- What should the dimensionality of Features be?



- Extract and Reconstruct *itself* \rightarrow *Auto*-Encoder + Decoder
- What should the dimensionality of Features be?
 - Smaller than Input! Why?
 - So encoding (feature extraction) process is a compression



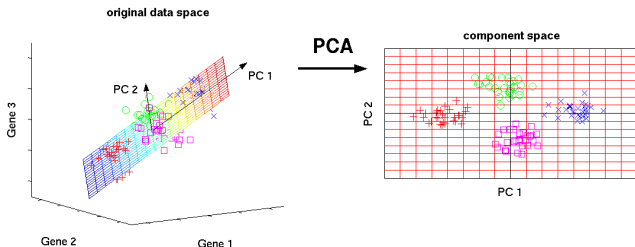
- Extract and Reconstruct *itself* \rightarrow *Auto*-Encoder + Decoder
- What should the dimensionality of Features be?
 - Smaller than Input! Why?
 - So encoding (feature extraction) process is a compression
- Compression: we need to learn prominent features (factors) only
 - Yet should be enough to reconstruct data



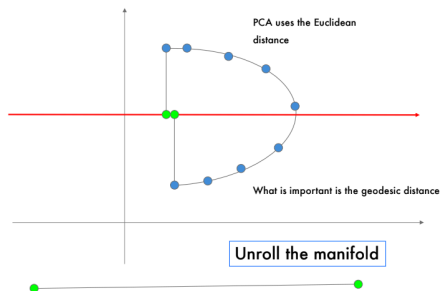
- Extract and Reconstruct *itself* \rightarrow *Auto*-Encoder + Decoder
- What should the dimensionality of Features be?
 - Smaller than Input! Why?
 - So encoding (feature extraction) process is a compression
- Compression: we need to learn prominent features (factors) only
 - Yet should be enough to reconstruct data
 - How? Unless for trivial cases, we need params θ

How to *best* extract features from data?

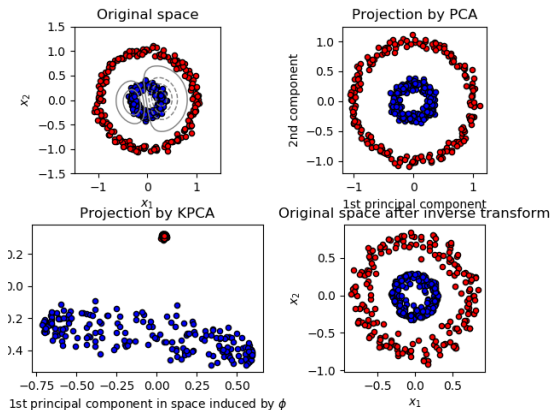
- ① Data Issue
- ② Methodology
- ③ Auto Encoder
- ④ Data Encoding to Embedded Representation
- ⑤ Decoding from Embedded Space



- Each Principal Component is a feature, and orthogonal to the rest
- Find a s.t. $x*a = \lambda a$ where λ is a scalar
 - One solution: perform SVD for x
 - Take only first few eigenvectors (hence approximation)



- Limitation: it's a linear method



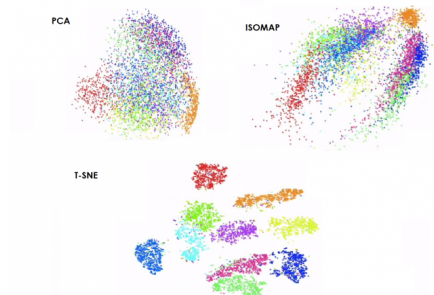
- Recall: why do we use a kernel?
- Still a PCA, and have to carefully choose the kernel

²https://en.wikipedia.org/wiki/Kernel_principal_component_analysis



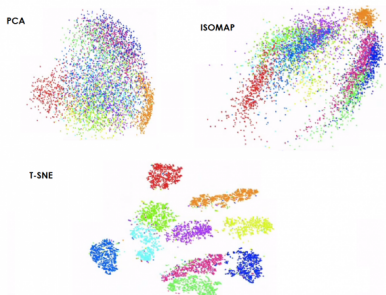
- Number of eigenvectors taken: 2000, 1000, 500, 100, 50, 10, 2, 1

Image credit: Gene Kogan



- Take into account relative distances between data points
- Is a non-linear method
- Both PCA and t-SNE suffer in high dimensional space.

³https://lvdmaaten.github.io/publications/papers/JMLR_2008.pdf



- Take into account relative distances between data points
- Is a non-linear method
- Both PCA and t-SNE suffer in high dimensional space.
- *It's time for Deep NN to be used!*

³https://lvdmaaten.github.io/publications/papers/JMLR_2008.pdf

- Original space:

- one-sided distance: $p_{i|j} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$

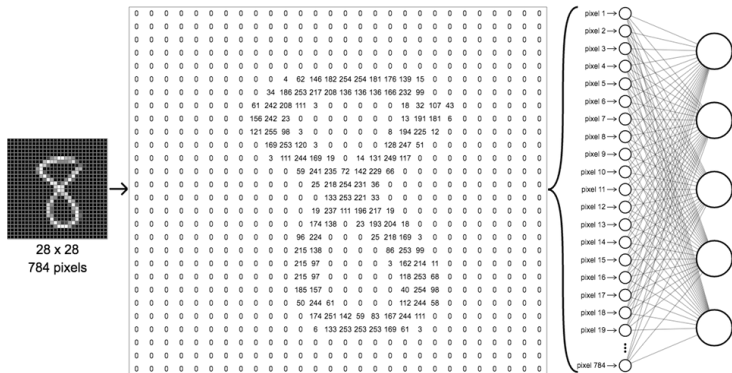
- two-sided (gaussian centered density): $p_{ij} = \frac{p_{i|j} + p_{j|i}}{2}$

- Embedded space: similar distance of 2 spaces:

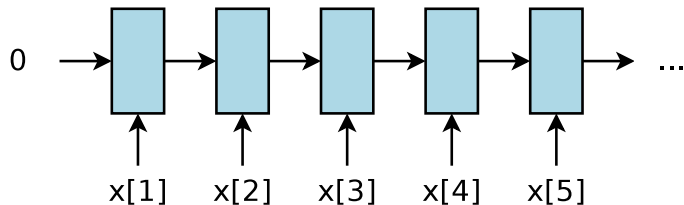
$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_i - y_k\|^2)^{-1}}$$

- Target: by minimizing the KL distance:

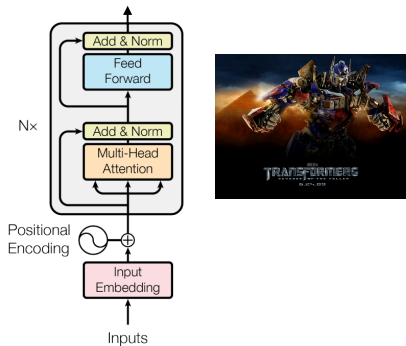
$$y_1, \dots, y_D = \operatorname{argmin} KL(P \| Q) = \operatorname{argmin} \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$



- Caveat: need to flatten the input
- With data knowledge, we can do a better job

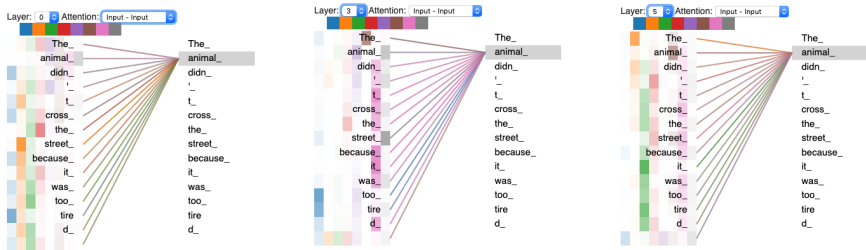


- Representation has sequential correlations

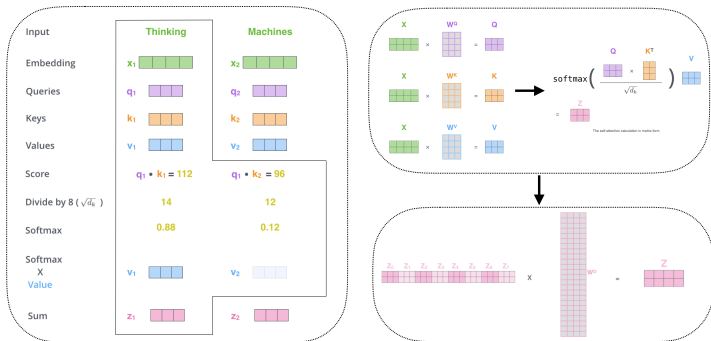


- We don't rely on RNN function anymore
- Instead, rely on N^* (multi-head attentions + FNN) blocks
- What's more ...

⁴<https://arxiv.org/pdf/1706.03762.pdf>

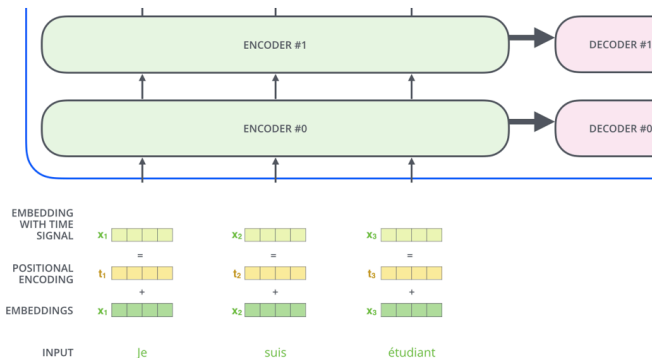


- Instead of encoding the sequence into a single hidden vector
- Now we calculate the relationship of each word vs. *every* another word
- Maybe a better idea for languages in which the order is very flexible

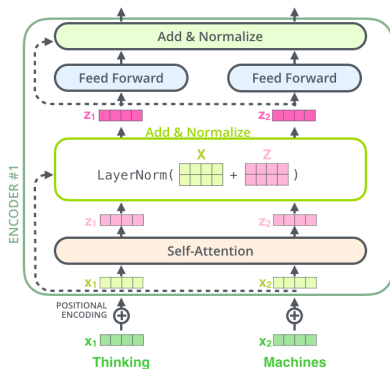


- Instead of using a single FC, use a triple (q, k, v) (query, key, value)
- Normalized with input length and with softmax
- Merge all head using concatenation and a FC with weight W^0

Multi-head Attn Able to Replace RNN? Carnegie Mellon

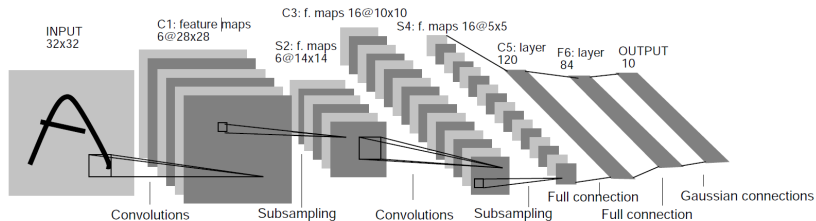


- Not yet, needs an alternative way to model sequence
- *Positional Encoding* is added to word embedding



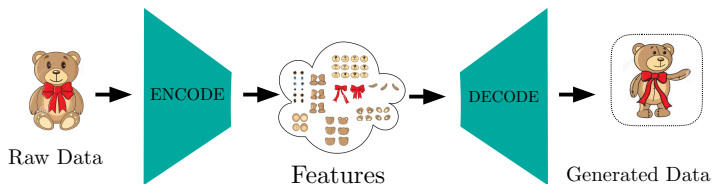
- Residual connection is used at self-attention and FNN
- Layer Norm⁵ is also applied.
- Note: it also stacks multiple encoders (and decoders (later))

⁵<https://arxiv.org/pdf/1607.06450.pdf>

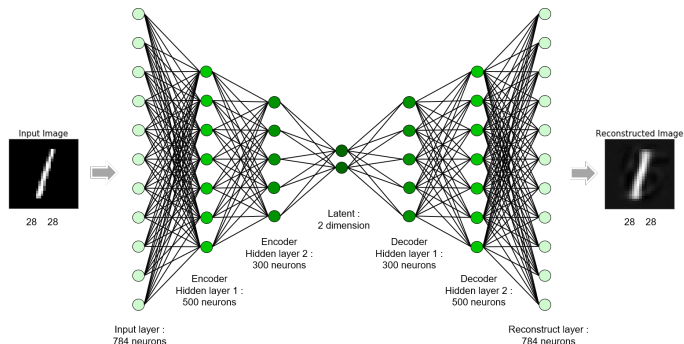


- Encoded representation has spatial correlations

- ① Data Issue
- ② Methodology
- ③ Auto Encoder
- ④ Data Encoding to Embedded Representation
- ⑤ Decoding from Embedded Space

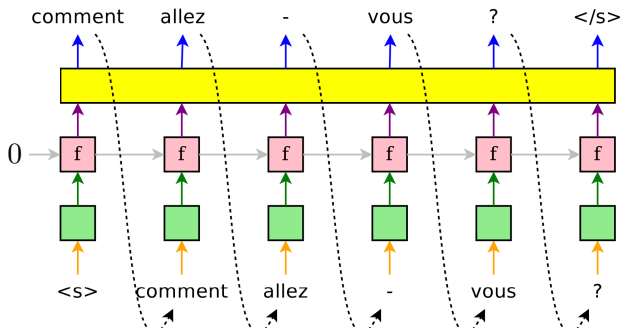


- Perform the inverse process of decoder
 - Encoder: compress from dimension N to D for $N > D$
 - Decoder: uncompress from D back to N

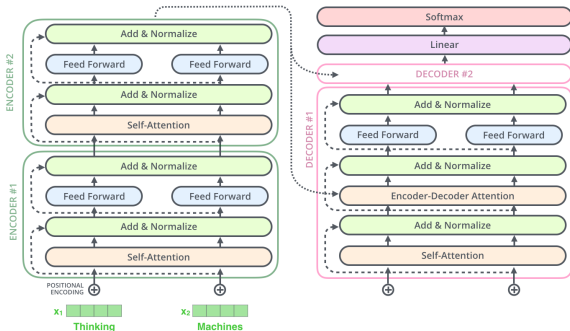


- Perform the inverse process of decoder
 - Encoder: MLP with decreasing dimensions
 - Decoder: MLP with increasing dimensions back to the original
- Tutorial: [▶ MNIST AutoEncoder Tensorflow-Keras](#)
- Tutorial: [▶ MNIST Denoising AutoEncoder Tensorflow-Keras](#)

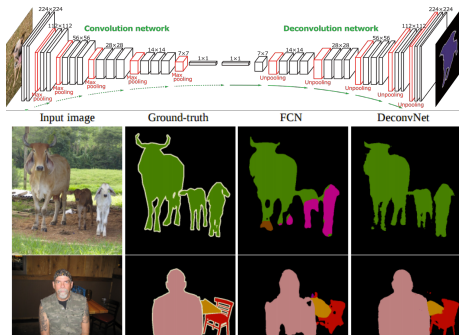
Image credit: i-systems.github.io



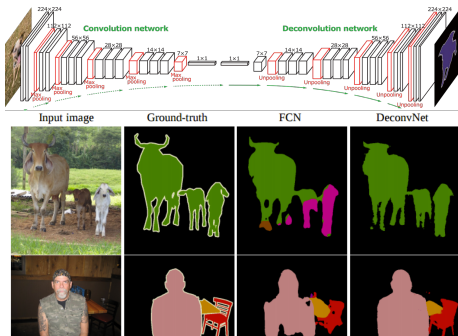
- Perform the inverse process of Encoder
 - Using previous output and previous hidden
 - It models $p(x_t|x_{<t})$ so sometimes it's called an *auto-regressive model*.
- Tutorial: [▶ RNN Text Generation with Tensorflow-Keras](#)



- Self-Attention is used, but only for the seen sequence to that point
- Then normal Attention is used as in Seq2Seq
- Decoding as usual to find a word in vocabulary.
- Tutorial: [Transformer NMT with Tensorflow-Keras](#)
- Tutorial: [NMT-Transformer with pre-trained weights](#)



- Deconvolution is only good enough for such tasks as segmetation where detail is not the first priority
- Maybe we need to add some constraint into embedded space?
- Or maybe we should not use MLE-based method?



- Maybe we need to add some constraint into embedded space?
- Or maybe we should not use MLE-based method?
- *Next lectures will help provide some appealing solutions to those questions.*