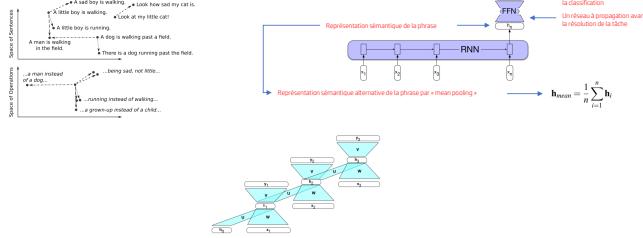


## Représentations neuronales au niveau de la phrase

### 1. Les réseaux récurrents

- Les réseaux de neurones récurrents (RNN) comme les LSTM permettent de traiter des séquences de toute longueur.
- Ils présentent le problème de « vanishing gradient » qui empêche de capturer pleinement la sémantique des phrases.



## Représentations neuronales au niveau de la phrase

### 2. Les transformateurs

- Contrairement aux RNN, les transformateurs ne présentent pas le problème de « vanishing gradient » et sont ainsi capables de représenter la sémantique des phrases plus complètement.
- Par contre, ils ne peuvent pas représenter des phrases de toute longueur puisqu'ils n'encodent pas la récurrence.

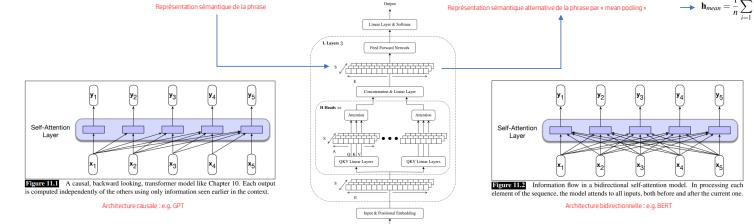


Figure 11.1 A causal, backward-looking, transformer model like Chapter 10. Each output is computed independently of the others using only information seen earlier in the context.

Architecture causale - e.g. GPT

Figure 11.2 Information flow in a bidirectional self-attention model. In processing each element of the sequence, the model attends to all inputs, both before and after the current one.

Architecture bidirectionnelle - e.g. BERT

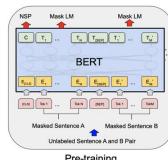
## Les modèles du langage pré-entraînés

### 1. Idée de base

- Plutôt que de travailler avec des modèles dont les poids sont initialisés aléatoirement et de ne se fonder que sur les plongements lexicaux comme base de connaissances, il est souvent préférable de pré-calculer les poids du réseau.
- Ainsi, les modèles pré-entraînés sont des modèles où des connaissances sur la langue ont déjà été apprises de façon non supervisée (autoencoder).
- Ces modèles pré-entraînés sont donc des modèles du langage.
- Suivant les modèles pré-entraînés, différentes techniques peuvent être utilisées.

### 2. Le modèle du langage masqué (« Masked Language Modeling ») (suite)

- L'idée est de prédire des tokens d'une phrase qui sont masqués aléatoirement, appelée « cloze task ».
- Please turn your homework \_\_\_\_ . Please turn \_\_\_\_ homework in.
- Cette technique est utilisée dans le cadre du modèle BERT (Devlin et al., 2019) basé sur des transformateurs bidirectionnels. En particulier, BERT est composé de 12 blocs chacun avec 12 têtes d'attention et la taille des couches cachées est de 768. Les phrases sont décomposées en tokens qui représentent des sous-mots (30000 tokens pour l'anglais).



## Les modèles du langage pré-entraînés

### 2. Le modèle du langage masqué (« Masked Language Modeling ») (suite)

- En plus des tokens masqués, certains mots sont remplacés par des mots non pertinents de façon aléatoire. Cela a pour objectif d'améliorer le regroupement sémantique des tokens (idée proche du « contrastive learning »).

### 3. Interprétation du modèle du langage masqué

- Etant donnée une séquence de tokens, les vecteurs de sortie de chaque token correspondent à des plongements contextualisés c'est-à-dire prenant en compte le contexte de la phrase.
- Ceci est particulièrement important pour les tokens polysémiques (e.g. jaguar) dont les représentations seront différentes selon le contexte de la phrase dans lesquels ils sont inclus.
- Ainsi, le vecteur yi du dernier bloc est le plongement du token xi.
- Il est possible de faire la moyenne des vecteurs yi des derniers blocs du transformeur pour atteindre une meilleure représentation.

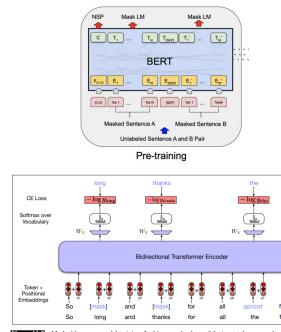


Figure 11.3 Mask language modeling training. In this example, three of the input words are selected (one of which is masked and the first is replaced with an unrelated word). The probabilities assigned by the model to these three items are used as the training loss. In the subsequent figures we deploy the input as words rather than word-level tokens; the reader should keep in mind that BERT and similar models actually use subword tokens instead.

## Les modèles du langage pré-entraînés

### 4. La prédiction de la phrase suivante

- Dans le cadre de BERT, une autre tâche est apprise qui consiste à prédire si une phrase en précède une autre ou non.
- En effet, cela peut aider à résoudre des tâches comme la détection de paraphrases, l'implication de phrases (« textual entailment ») ou la cohérence textuelle.
- Pour se faire, un token de segment est créé qui facilite la différenciation des deux phrases: [SEP].
- Notez que la prédiction n'est faite que sur le vecteur du token [CLS] qui correspond à la représentation de la phrase dans sa globalité.
- Chaque token est représenté par un plongement lexical, un plongement positionnel et un plongement de segment.

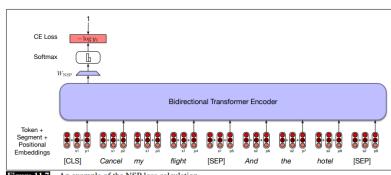


Figure 11.7 An example of the NSP loss calculation.

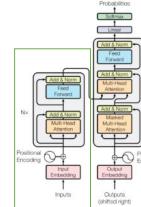
## Les modèles du langage pré-entraînés

### 5. Types de modèles pré-entraînés

- Il existe différents types de modèles du langage pré-entraînés, selon qu'ils ont été entraînés sur des tâches de compréhension de la langue (« natural language understanding ») ou des tâches de génération (« natural language generation ») ou les deux à la fois.

### 6. Modèles de type Encodateur (« natural language understanding »)

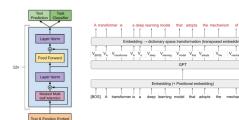
- BERT (Devlin et al., 2019)
- DistillBERT (Sanh et al., 2020) – un modèle compact de BERT par distillation
- RoBERTa (Liu et al., 2019) – plus de données, pas de NSP et des phrases longues
- XLM (Lample et Conneau, 2019) – version multilingue avec une nouvelle fonction d'erreur
- ALBERT (Lan et al., 2019) – prédiction de l'ordre des phrases
- ELECTRA (Clark et al., 2020) – architecture de type Generative Adversarial Network (GAN)
- DeBERTa (He et al., 2021) – deux vecteurs de représentation des tokens



## Les modèles du langage pré-entraînés

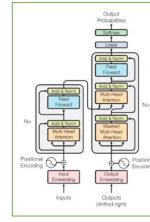
### 7. Modèles de type Décodeur (« natural language generation »)

- GPT (Rafford et al., 2018) – modèle causal
- CTRL (Kesher et al., 2019) – contrôle du style de génération
- GPT-2 (Salimian et al., 2019) – peut produire des textes cohérents
- GPT-3 (Brown et al., 2020) – version avec beaucoup plus de paramètres



### 8. Modèles de type Encodateur – Décodeur (NLU + NLG)

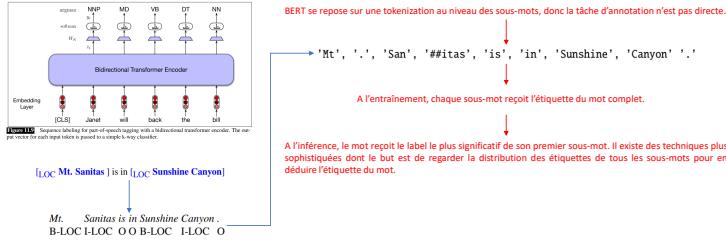
- T5 (Rafford et al., 2020)
- BART (Lewis et al., 2019) – BERT + GPT
- M2M-100 (Fan et al., 2020) – modèle multilingue
- BigBIRD (Zaheer et al., 2020) – optimisé en terme de mémoire utilisée



## Le « fine tuning » ou ajustement

### 3. Annotation par « fine tuning »

- Dans le cadre d'une annotation au niveau des mots ou des séquences selon le modèle BIO, le « fine tuning » est facile à modéliser avec une particularité en ce qui concerne la tokenization.



## Le « fine tuning » ou ajustement

### 4. Classification de séquences par « fine tuning »

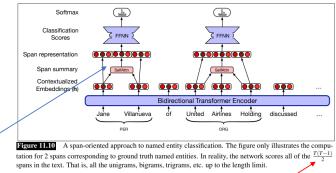
- Cette tâche consiste à classer des séquences de tokens d'intérêt, comme par exemple les unités polylexiques, des unités étiquetées par des experts (e.g. mémoire épisodique, annotations de psychiatrie, etc.). Cette tâche se rapproche des méthodes d'annotation BIO.

Une séquence est représentée par la concaténation de son plongement initial, de son plongement final et d'un plongement calculé par un mécanisme d'auto-attention sur l'ensemble de la séquence.

$$\text{spanRep}_{ij} = [h_i; h_j; g_{ij}]$$

$$g_{ij} = \text{SelfAttention}(h_{ij})$$

Ainsi, le plongement intermédiaire entre le plongement initial et le plongement final est une moyenne (« average pooling ») de tous les plongements de la séquence pondérés par le mécanisme d'auto-attention.



Toutes les sous-séquences jusqu'à une taille maximale sont classifiées

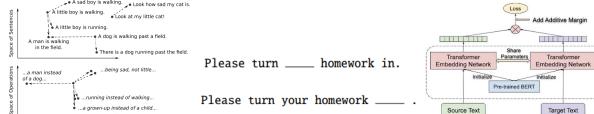
## Amélioration de la représentation phrasique

### 1. Importance de la représentation sémantique

- Les modèles du langage ont comme principale tâche de correctement représenter les phrases dans un espace sémantique cohérent. Plus cet espace sera juste et plus les tâches pourront être apprises simplement et plus les performances seront élevées.

### 2. Limites des modèles du langage

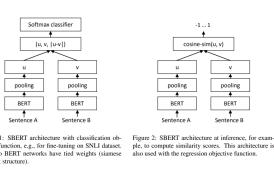
- Les modèles du langage ont été appris sur la base de tâches bien définies. Pour BERT, le « masked language modelling » et pour GPT le « next word prediction ».
- Or ces tâches ne peuvent garantir la construction d'un espace sémantique cohérent même si elles s'en rapprochent.
- Des travaux s'attachent à construire des espaces de représentation sémantique, notamment basés sur l'idée de « contrastive learning » (Reimers and Gurevych 2019, Chen et al. 2020)



## Amélioration de la représentation phrasique

### 1. SentenceBERT (Reimers & Gurevych, 2019)

- SentenceBERT propose d'ajuster un modèle du langage de type BERT en se basant sur une architecture siamoise.
- Pour se faire, l'architecture est entraînée sur les jeux de données SNLI (<https://cims.nyu.edu/~sbowman/multinli/>) et MultiNLI (<https://cims.nyu.edu/~sbowman/multinli/>).
- Cette architecture est ensuite testée sans nouvel apprentissage (i.e. « linear probing ») sur la série de jeux de données STS qui évaluent les similarités entre phrases .



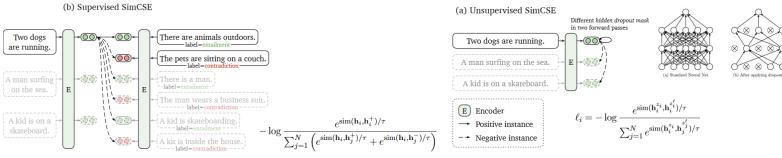
	STS2	STS3	STS4	STS5	STS6	STS7	SICK-R	Avg.
SGNS embeddings	38.76	57.98	57.94	63.15	61.04	46.35	58.40	54.81
Avg. BERT embeddings	38.76	57.98	57.94	63.15	61.04	46.35	58.40	54.81
BERT CLS vector	32.46	66.73	62.15	56.26	66.07	48.03	65.63	63.01
Universal Sentence Encoder	64.49	67.80	64.63	76.93	73.14	74.97	76.69	71.23
SICK-LU Large	72.27	78.46	74.90	80.99	76.25	79.23	75.77	76.55
SGBERT-NL	74.53	77.00	73.14	81.85	76.86	79.10	74.29	76.68

Table 1: SentenceBERT achieves the best results on the STS benchmark for all the classification objectives, on the SICK benchmark for the classification objective function, and on the STS benchmark for the correlation objective function. The table also compares SentenceBERT with other state-of-the-art models on the STS2012-2016, STSb, STSbenchmark, SICK-R, SICK relatedness datasets. SentenceBERT outperforms all the other models on all the classification objectives except for the correlation objective function.

## Amélioration de la représentation phrasique

### 1. SimSCE (Gao et al., 2022)

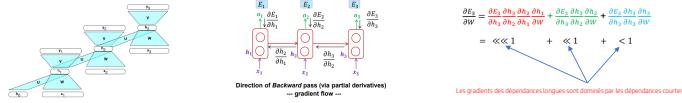
- SimSCE propose d'ajuster un modèle du langage de type BERT ou RoBERTa en se basant sur deux idées différentes: une architecture non supervisée et une architecture supervisée dans le cadre d'un apprentissage contrastif.
- L'architecture non supervisée reçoit une même phrase deux fois en entrée dans le même batch mais avec deux masques de « dropout » différents. L'objectif d'apprentissage est donc de rapprocher ses deux représentations.
- L'architecture supervisée est entraînée sur les jeux de données SNLI (<https://nlp.stanford.edu/projects/snli/>) et MultiNLI (<https://cims.nyu.edu/~sbbowman/multinli/>) pour donner les meilleurs résultats.
- Dans l'architecture supervisée, des exemples positifs (« entailment ») et des exemples négatifs (« contradiction ») partagent le même « batch » et les exemples positifs doivent être rapprochés alors que les négatifs éloignés.



## Les limitations des représentations neuronales

### 1. Représentations des textes par RNN

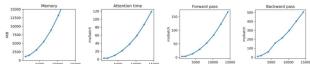
- Les RNN ont une capacité limitée de représentation du fait de leur récurrence. Tout texte peut être représenté par un RNN. Par contre, leur capacité à représenter les séquences longues est limitée par la notion de « vanishing gradient ».



### 2. Représentation des textes par transformeurs

- Les transformateurs ont une capacité limitée de représentation du fait de leur non récurrence. Par exemple, pour BERT le nombre de tokens en entrée est limité à 512.
- Cette limitation est aussi due à la complexité quadratique en coût de calcul:  $O(N^2)$  où  $N$  est la taille de la séquence.

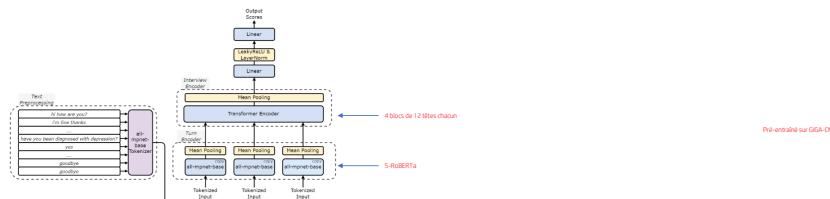
Pour un transformeur à une tête et un set bloc sur une carte graphique RTX8000 GPU.



## Les modèles hiérarchiques

### 3. Une représentation à partir de modèles pré-entraînés (Milintsevich et al., 2023)

- Afin de prendre en compte les modèles pré-entraînés de type BERT comme représentation des phrases et les transformateurs comme agrégateurs de sémantique, des architectures hiérarchiques plus performantes peuvent être élaborées.
- Dans ce cas, les phrases sont représentées grâce à des modèles de type SentenceBERT, ici S-RoBERTa qui sont ensuite aggrégées grâce à un transformeur.



## Analyse au niveau du document

### 1. Différentes applications

- Au semestre 1 du Master 1, nous avons étudié le langage au niveau du mot (analyse lexicale).
- Au semestre 2 du Master 1, nous avons étudié le langage au niveau de la phrase (analyse phrasique)
- Or, de nombreuses applications du langage naturel traitent des textes qui sont des séquences longues de phrases ou de paragraphes.

### 2. Applications intra-document

- Cohérence textuelle : analyse des coréférences
- Structure du discours : chaînes lexicales, segmentation thématique
- Structure argumentative : « argument mining »

### 3. Applications de classification

- Analyse de sentiments, d'émotions, de fake news, de diagnostic médical, de profilage etc ...

### 4. Applications de génération

- Traduction automatique
- Résumé de textes
- Dialogues

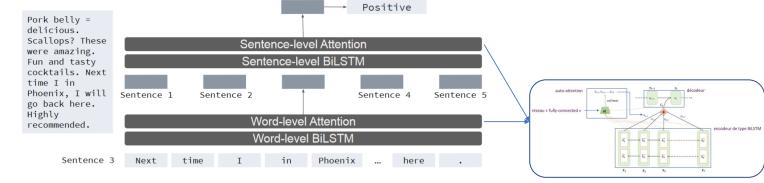
## Les modèles hiérarchiques

### 1. Utiliser la nature hiérarchique des textes

- Les textes sont composés hiérarchiquement de caractères, de mots, de phrases et de paragraphes.
- Les modèles de représentation prennent en compte cette organisation pour modéliser de longs textes.

### 2. Une représentation à partir de BiLSTM (Yang et al., 2016)

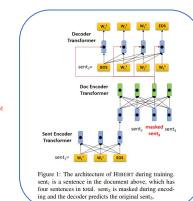
- Le texte est d'abord divisé en phrases qui sont elles-mêmes divisées en mots. Chaque phrase est représentée par la dernière couche du BiLSTM initialisé à partir de plongements lexicaux. Toutes les phrases sont ensuite aggrégées dans un BiLSTM dont la dernière couche représente le texte dans son intégralité.



## Les modèles hiérarchiques

### 4. Des modèles hiérarchiques pré-entraînés au niveau du document

- (Zhang et al., 2020) proposent de pré-entraîner un modèle hiérarchique à partir de la notion de « masked language modeling » mais au niveau de la phrase.
- L'idée est de construire un modèle capable de représenter sémantiquement un texte de manière auto-supervisée, c'est-à-dire en reconstruisant des phrases masquées et en prédisant la phrase suivante.
- Le modèle HIBERT a été entraîné à partir du corpus GIGA-CM comportant 6.3 millions de documents et 2.8 milliards de tokens.



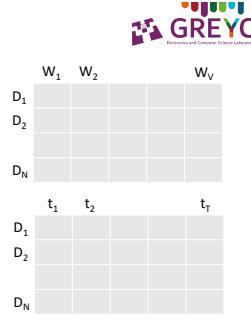
Model	B	S	E
Pointer-Genegate	96.5	77.25	76.55
Abstrac-MR+BL	98.9	13.82	38.90
DCA	41.60	19.47	37.92
ScorePredictor	40.80	19.47	37.92
Inconsistency	40.60	19.47	37.92
Insistence-Up	41.20	19.47	37.92
Lead3	40.31	17.70	36.57
SummarizeNer	39.60	18.20	35.30
Necklace	40.30	18.20	35.30
Refresh	40.00	18.20	36.60
Necklace-MMR	41.20	18.20	36.60
Baseline	41.50	18.70	37.60
HCS	41.70	18.50	37.90
Learned	41.00	18.70	37.25
HierTransferer	41.1	18.69	37.35
HIBERT	41.2	18.70	37.30
HIBERT (on-domain)	42.0	19.70	38.55
HIBERT (cross-domain)	42.3	19.95	38.43

Table 1: Results of various models on the CNDM test set using full length. F1 score (F1), R@10 (R), and P@10 (P).



## Latent Dirichlet Allocation

- Dimensionality reduction:
  - # documents: N
  - # words in vocabulary: V
- documents can be represented using document-term matrix  $\mathbb{R}^{N \times V}$



Assume T topics are learned from this data.

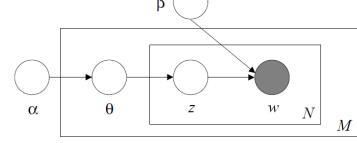
Documents are represented as distribution over topics  $\mathbb{R}^{N \times T}$

- Unsupervised learning: can be compared to clustering.
- Words are clustered together to form topics based on their co-occurrence patterns
- Documents are clustered based on their topic distributions.

Blei et al., JMLR 2003

## Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a generative probabilistic model of a corpus. Within LDA, documents are represented as random mixtures over latent topics, where each topic is characterized by a distribution over words.



Lets assume we want to generate M documents:

- Choose N: number of words in the document
- Choose  $\theta \sim \text{Dir}(a)$ : topic proportions for document w
- For each of the N words:
  - Choose a topic  $z_n \sim \text{Multinomial}(\theta)$ : topic assignment for document w
  - Choose a word  $w_n$  from  $p(w|z_n, \beta)$

## Latent Dirichlet Allocation

Given  $\alpha$  and  $\beta$ , the joint distribution of a topic mixture  $\theta$ , set of  $N$  topics  $z$ , and  $N$  words  $w$  is given by:

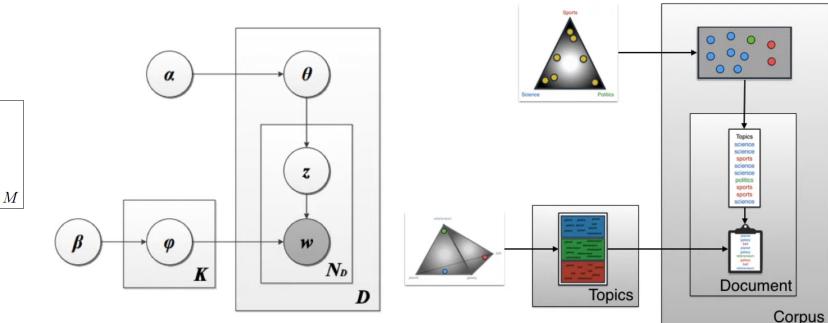
$$p(\theta, z, w | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

Integrating over  $\theta$  and summing over  $z$ , we obtain the marginal distribution of a document:

$$p(w | \alpha, \beta) = \int p(\theta | \alpha) \left( \prod_{n=1}^N \sum_z p(z_n | \theta) p(w_n | z_n, \beta) \right) d\theta.$$

Finally, taking product of marginal probabilities of single documents, we obtain the probability of a corpus:

$$p(D | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left( \prod_{n=1}^{N_d} \sum_z p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) d\theta_d.$$



## Latent Dirichlet Allocation

LDA is part of a larger field of *probabilistic modeling*.

We treat the data as arising from a generative process that includes *hidden variables*.

This generative process defines a joint probability distribution over both the observed and hidden variables

We perform data analysis by using joint distribution to compute conditional distribution of the hidden variables given the observed variables.

This conditional distribution is called *posterior distribution*.

observed variables:  $w$

hidden variables:  $\theta, z, \beta$

$$p(\theta, z, \beta | w) = \frac{p(\theta, z, \beta, w)}{p(w)}$$

## Latent Dirichlet Allocation

- The posterior cannot be computed because the denominator is intractable.
- Topic modeling algorithms form an approximation of the equation by adapting an alternative distribution over the latent topic structures to be close to the true posterior.
- Topic modeling algorithms generally fall into two categories:
  - Sampling based algorithms
  - Variational algorithms
- The most commonly used sampling algorithm for topic modeling is *Gibbs Sampling*

## Gibbs Sampling

- Gibbs sampling procedure considers each word token in the text collection in turn.
- Estimate the probability of assigning the current word token to each topic, conditioned on the topic assignments to all the other word tokens.

$$P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot) \propto \frac{C_{w_i,j}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_i,j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i,t}^{DT} + T\alpha}$$

- From this conditional distribution, a topic is sampled and stored as the new assignment for this word token.

## Gibbs Sampling

- $n_{dk}$ : number of words assigned to topic k in document d
- $n_{kw}$ : number of times word w is assigned to topic k
- $n_k$ : number of times any word is assigned to topic k

```

Input: words w ∈ documents d
Output: topic assignments z and counts n_{d,k}, n_{k,w}, and n_k
begin
    randomly initialize z and increment counters
    foreach iteration do
        for i = 0 → N - 1 do
            word ← w[i]
            topic ← z[i]
            n_{d,topic}+=1; n_{word,topic}+=1; n_{topic}+=1
            for k = 0 → K - 1 do
                p(z = k|·) = (n_{d,k} + α_k)^(n_{k,w} + β_w) / (n_k + β × W)
            end
            topic ← sample from p(z|·)
            z[i] ← topic
            n_{d,topic}+=1; n_{word,topic}+=1; n_{topic}+=1
        end
    end
    return z, n_{d,k}, n_{k,w}, n_k
end

```

Griffiths and Steyvers (2004)

## Gibbs Sampling

### Example

Assume we have some document with random word-topic assignment

India	enters	world	cup	final	
1	3	1	2	4	

We have count matrix  $C^{WT}$

	1	2	3	4	
India	70	5	0	8	
enters	2	3	15	6	
world	28	4	12	1	
cup	6	43	6	0	
final	7	0	9	31	

## Gibbs Sampling

### Example

Assume we have some document with random word-topic assignment

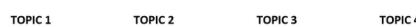
India	enters	world	cup	final	
1	3	1	2	4	

We have count matrix  $C^{WT}$

	1	2	3	4	
India	70	5	0	8	
enters	2	3	15	6	
world	28	4	12	1	
cup	6	43	6	0	
final	7	0	9	31	

## Gibbs Sampling

- Consider the contribution of each topic towards this document

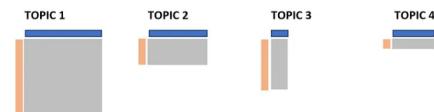


- Next, we take how many times each topic is assigned to this word



## Gibbs Sampling

- Multiply these values to get conditional probabilities



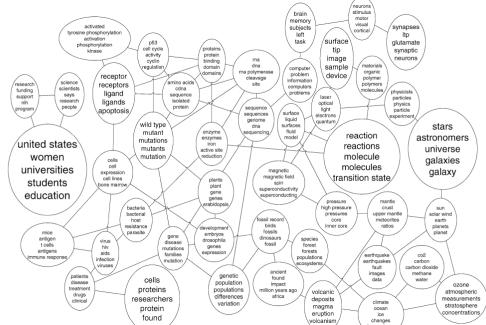
- Finally, pick one of the topics from this distribution and update the variables accordingly.
- Repeat this for every word.



## Topic model variations

### 1. Correlated topic model

- Topic graph learned from 16,351 OCR articles from science (1990-1999)



Blei, David M., and John D. Lafferty. "A correlated topic model of science." (2007): 17-35.

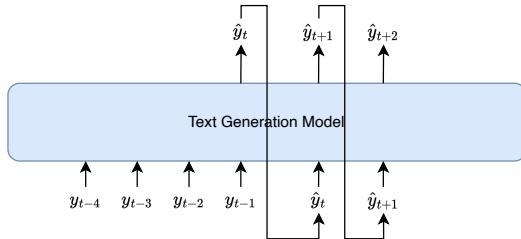
## Topic model variations

### 1. Structured Topic Model (STM)

- LDA learns only based on input text.
- Certain sources may be more likely to write about politics.
- Metadata can include date published, author, publication, likes on social media, etc.
- Within LDA our topic distribution comes from Dirichlet distribution.
- STM defines topic distributions based on document metadata
- We need to go from  $X$ ,  $1xp$  metadata vector to  $1xk$  vector of topic distribution.
- We multiply  $X$  with  $pxk$  weight matrix  $t$ .

## Autoregressive Models

- In autoregressive text generation models, at each time step  $t$ , our model takes in a sequence of tokens of text as input  $y_{<t}$  and outputs a new token,  $\hat{y}_t$



## During a Single Step

- At each time step  $t$ , our model computes a vector of scores for each token in our vocabulary,  $S \in \mathbb{R}^V$ :

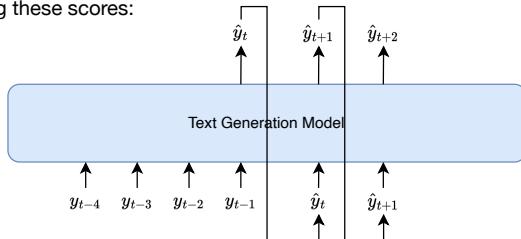
$$S = f(\{y_{<t}\}, \theta)$$

- Then, we compute a probability distribution  $P$  over  $w \in V$  using these scores:

$$P(y_t | \{y_{<t}\}) = \frac{\exp(S_w)}{\sum_{w' \in V} \exp(S_{w'})}$$

## Autoregressive Models

- At each time step  $t$ , our model computes a vector of scores for each token in our vocabulary,  $S \in \mathbb{R}^V$ . Then, we compute a probability distribution  $P$  over  $w \in V$  using these scores:



## Inference and Training

- At inference time, our decoding algorithm defines a function to select a token from this distribution:

$$\hat{y}_t = g(P(y_t | \{y_{<t}\}))$$

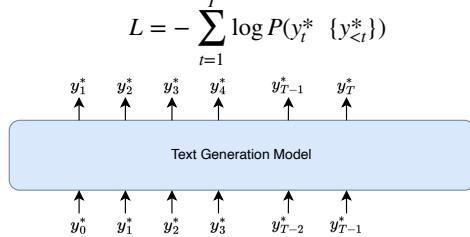
- We train the model to minimize the negative loglikelihood of predicting the next token in the sequence:

$$L_t = -\log P(y_t^* | \{y_{<t}^*\})$$

## Maximum Likelihood Training

### Teacher Forcing

- Trained to generate the next word  $y_t^*$  given a set of preceding words  $\{y_i^*\}_{i < t}$



## Decoding

- At each time step  $t$ , our model computes a vector of scores for each token in our vocabulary,  $S \in \mathbb{R}^V$ :

$$S = f(\{y_{<t}\}, \theta)$$

- Then, we compute a probability distribution  $P$  over  $w \in V$  using these scores:

$$P(y_t | \{y_{<t}\}) = \frac{\exp(S_w)}{\sum_{w' \in V} \exp(S_{w'})}$$

- Our decoding algorithm defines a function to select a token from this distribution:

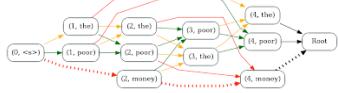
$$\hat{y}_t = g(P(y_t | \{y_{<t}\}))$$

## Greedy Methods

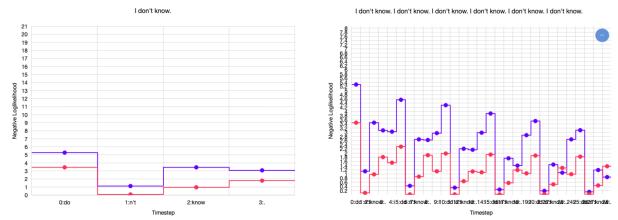
- Argmax Decoding

- Selects the highest probability token in  $P(y_t | \{y_{<t}\})$

- Beam Search



## Repetition in Greedy Methods

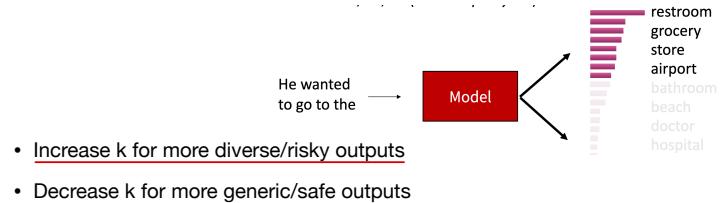


## Top-k sampling

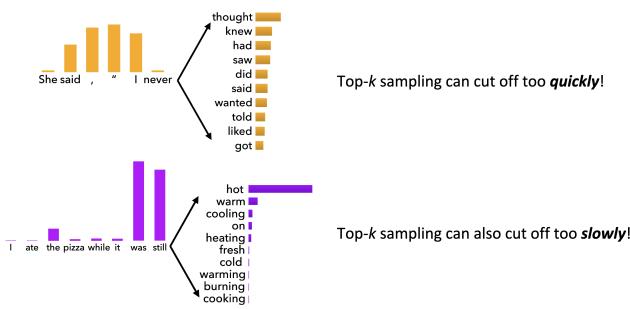
- Problem: Vanilla sampling makes every token in the vocabulary an option
- Even if most of the probability mass in the distribution is over a limited set of options, the tail of the distribution could be very long
- Many tokens are probably irrelevant in the current context
- Why are we giving them individually a tiny chance to be selected?
- Why are we giving them as a group a high chance to be selected?
- Solution: Top-k sampling
- Only sample from the top k tokens in the probability distribution

## Top-k sampling

- Only sample from the top k tokens in the probability distribution
- Common values are  $k = 5, 10, 20$  (but it's up to you!)



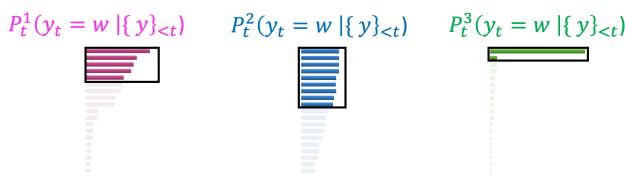
## Top-k sampling



## Top-p (nucleus) sampling

- Problem: The probability distributions we sample from are dynamic
  - When the distribution  $P_t$  is flatter, a limited  $k$  removes many viable options
  - When the distribution  $P_t$  is peakier, a high  $k$  allows for too many options to have a chance of being selected
- Solution: Top-p sampling
  - Sample from all tokens in the top  $p$  cumulative probability mass (i.e., where mass is concentrated)
  - Varies  $k$  depending on the uniformity of  $P_t$

## Top-p (nucleus) sampling



## Softmax Temperature

- On timestep  $t$ , the model computes a prob distribution  $P_t$  by applying the softmax function to a vector of scores  $s \in \mathbb{R}^V$

$$P(y_t | \{y_{<t}\}) = \frac{\exp(S_w)}{\sum_{w' \in V} \exp(S_{w'})}$$

- You can apply a temperature hyperparameter  $\tau$  to the softmax to rebalance  $P_t$ :

$$P(y_t | \{y_{<t}\}) = \frac{\exp(S_w / \tau)}{\sum_{w' \in V} \exp(S_{w'} / \tau)}$$

## Softmax Temperature

- Raise the temperature  $\tau > 1$ :  $P_t$  becomes more uniform
  - More diverse output (probability is spread around vocab)
- Lower the temperature  $\tau < 1$ :  $P_t$  becomes more spiky
  - Less diverse output (probability is concentrated on top words)

## Re-ranking

- Decode a bunch of sequences
  - 10 candidates is a common number
- Define a score to approximate quality of sequences and re-rank by this score
  - Simplest is to use perplexity
    - Careful! Remember that repetitive methods can generally get high perplexity.

## Decoding

- Decoding is still a challenging problem in natural language generation
- Human language distribution is noisy and doesn't reflect simple properties (i.e., probability maximization)
- Different decoding algorithms can allow us to inject biases that encourage different properties of coherent natural language generation
- Some of the most impactful advances in NLG of the last few years have come from simple, but effective, modifications to decoding algorithms

## Content Overlap Metrics

Ref: They walked to the grocery store .  
Gen: The woman went to the hardware store .

- Compute a score that indicates the similarity between generated and gold-standard (human-written) text
- Fast and efficient and widely used
- Two broad categories:
  - N-gram overlap metrics (e.g., BLEU, ROUGE, METEOR, CIDEr, etc.)
  - Semantic overlap metrics (e.g., PYRAMID, SPICE, SPIDEr, etc.)

## Content Overlap Metrics

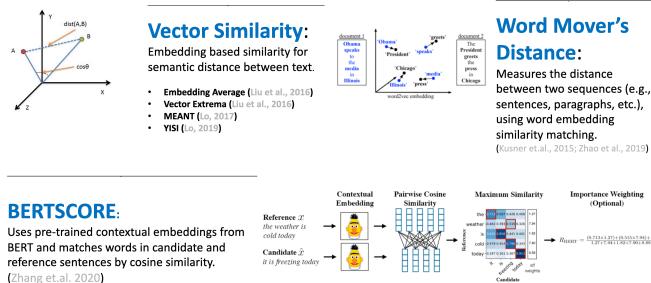
- They're not ideal for machine translation
- They get progressively much worse for tasks that are more open-ended than machine translation
  - Worse for summarization, where extractive methods that copy from documents are preferred
  - Much worse for dialogue, which is more open-ended than summarization
  - Much, much worse story generation, which is also open-ended, but whose sequence length can make it seem you're getting decent scores!

Dans le cas de la summarisation, où des méthodes extractives (copier des parties du texte source) sont souvent préférées, les métriques de recouvrement de contenu peuvent ne pas bien refléter la qualité du résumé, car la redondance peut être considérée comme une bonne pratique.

Les métriques de recouvrement de contenu peuvent être encore moins adaptées pour évaluer la qualité des dialogues. Les conversations sont souvent évaluées en fonction de la cohérence, de la pertinence contextuelle, et de la manière dont elles maintiennent une interaction significative, des aspects qui ne sont pas capturés efficacement par la simple comparaison du contenu.

La génération d'histoires est une tâche très ouverte où la créativité et la cohérence narrative sont cruciales. Les métriques de recouvrement de contenu peuvent ne pas saisir ces aspects, en particulier lorsque les histoires sont longues et complexes.

## Model-Based Metrics



## Human Evaluation

- Ask humans to evaluate the quality of generated text
- Human judgments are regarded as the gold standard
- Humans are inconsistent

# Non-autoregressive Models

Application	Example	Use seq2seq
Machine translation	S: Turing studied at King's College, where he was awarded first-class honours in mathematics. T: Turing studierte am King's College, wo er erste Klasse Auszeichnungen in Mathematik erhielt.	✓
Summarization	S: Court members Deborah Portz and Peter Verniero did not participate in the Nelson case. T: Court members didn't participate in the case.	?
Sentence fusion	S: Turing was born in 1912. Turing died in 1954. T: Turing was born in 1912 and he died in 1954.	✗
Grammar correction	S: New Zealand have a cool weather. T: New Zealand has cool weather.	✗

# Non-autoregressive Models

- Most NLP tasks apart from Machine Translation are monolingual
- Sources and targets often overlap
- Generating the target from scratch is wasteful
- Can reconstruct most of the target from the source via basic operations like KEEP, DELETE, INSERT

Turing was born in 1912 . Turing died in 1954 .  
 KEEP KEEP KEEP KEEP KEEP DEL INS PRON KEEP KEEP KEEP KEEP  
 Turing was born in 1912 and he died in 1954 .

# Applications

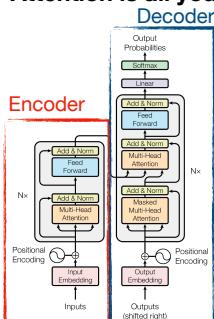
- Grammatical Error Correction
- Text Simplification
- Sentence fusion
- Style transfer
- Text normalisation
- Text summarisation
- Automatic post-editing for machine translation

# Advantages

- Text editing models need less training data
- They are faster at the inference
- They are more faithful
- We can control what model adds or removes
- We can incorporate external knowledge

# Transformer Architecture

Vaswani, Ashish, et al. "Attention is all you need." (2017)



# Byte-Pair Encoding

## Training

- First, we compute the unique set of words in the corpus
- Split each word into characters
- Start merging most frequent pairs
- ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)
- ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

# Byte-Pair Encoding

## Training

- First, we compute the unique set of words in the corpus
- Split each word into characters
- Start merging most frequent pairs
- The most frequent pair is ("u", "g")
- ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)
- ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

Vocabulary: ["b", "g", "h", "n", "p", "s", "u", "ug"]

Corpus: ("h" "ug", 10), ("p" "ug", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "ug" "s", 5)

# Byte-Pair Encoding

## Training

- First, we compute the unique set of words in the corpus
- Split each word into characters
- Start merging most frequent pairs
- The most frequent pair is ("u", "n")
- ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)
- ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

Vocabulary: ["b", "g", "h", "n", "p", "s", "u", "ug", "un"]

Corpus: ("h" "ug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("h" "ug" "s", 5)

# Byte-Pair Encoding

## Training

- First, we compute the unique set of words in the corpus
- Split each word into characters
- Start merging most frequent pairs
- The most frequent pair is ("h", "ug")
- ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)
- ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

Vocabulary: ["b", "g", "h", "n", "p", "s", "u", "ug", "un", "hug"]

Corpus: ("hug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("hug" "s", 5)

# Byte-Pair Encoding

## Tokenization

1. Normalization
2. Pre-tokenization
3. Splitting the words into individual characters
4. Applying the merge rules learned in order on those splits

# Byte-Pair Encoding

## Tokenization

1. Normalization
2. Pre-tokenization
3. Splitting the words into individual characters
4. Applying the merge rules learned in order on those splits

### Learned rules:

("u", "g") -> "ug"  
("u", "n") -> "un"  
("h", "ug") -> "hug"

# Byte-Pair Encoding

## Tokenization

1. Normalization
2. Pre-tokenization
3. Splitting the words into individual characters
4. Applying the merge rules learned in order on those splits

### Learned rules:

("u", "g") -> "ug"  
("u", "n") -> "un"  
("h", "ug") -> "hug"

"bug" -> "b", "ug" | "mug" -> "[UNK]", "ug"

# WordPiece Tokenization

## Training

- Same as BPE but different merging rule

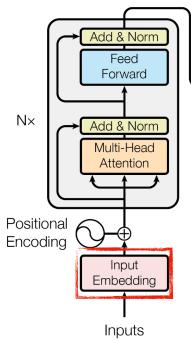
```
score=(freq_of_pair)/
(freq_of_first_element×freq_of_second_element)
```

# WordPiece Tokenization

## Tokenization

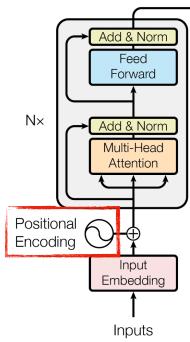
- Same as BPE but the WordPiece starts by searching for the longest sub-token in the vocabulary

# Transformer Encoder



- Byte-pair Encoding (GPT, GPT-2, RoBERTa, BART, and DeBERTa) or WordPiece (BERT) tokenization
- Randomly initialized learnable embedding layer

# Transformer Encoder



- Byte-pair Encoding (GPT, GPT-2, RoBERTa, BART, and DeBERTa) or WordPiece (BERT) tokenization
- Randomly initialized learnable embedding layer
- Positional encoding to add information about words' position in the sequence

# Positional Encoding

- Since Transformer doesn't have recursion it lacks information about words' positions

Let  $t$  be the desired position in an input sentence,  $\vec{p}_t \in \mathbb{R}^d$  be its corresponding encoding, and  $d$  be the encoding dimension (where  $d \equiv 2$ )  
Then  $f : \mathbb{N} \rightarrow \mathbb{R}^d$  will be the function that produces the output vector  $\vec{p}_t$  and it is defined as follows:

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

where

$$\omega_k = \frac{1}{10000^{2k/d}}$$

# Positional Encoding

## Intuition

0 :	0 0 0 0	8 :	1 0 0 0
1 :	0 0 0 1	9 :	1 0 0 1
2 :	0 0 1 0	10 :	1 0 1 0
3 :	0 0 1 1	11 :	1 0 1 1
4 :	0 1 0 0	12 :	1 1 0 0
5 :	0 1 0 1	13 :	1 1 0 1
6 :	0 1 1 0	14 :	1 1 1 0
7 :	0 1 1 1	15 :	1 1 1 1

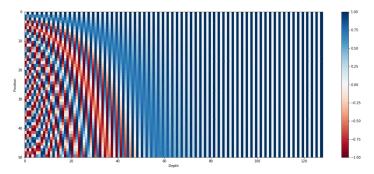
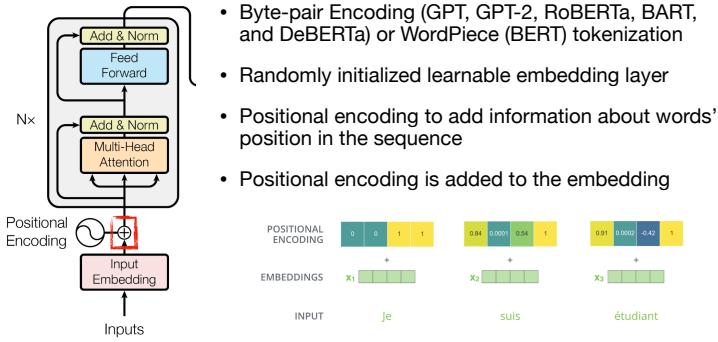
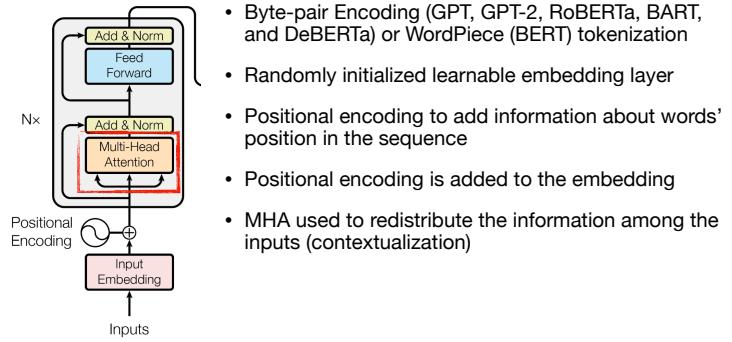


Figure 2 - The 128-dimensional positional encoding for a sentence with the maximum length of 15. Each row represents the embedding vector  $\vec{p}_t$ .

## Transformer Encoder



## Transformer Encoder



## Multi-Head Self-Attention

Vaswani, Ashish, et al. "Attention is all you need." (2017)

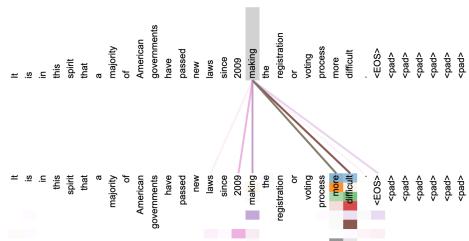
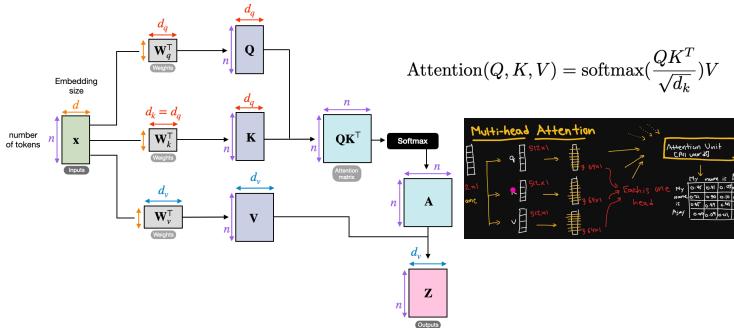


Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb 'making', completing the phrase 'making...more difficult'. Atentions here shown only for the word 'making'. Different colors represent different heads. Best viewed in color.

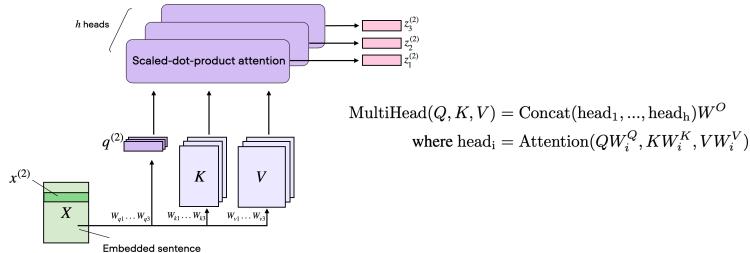
## Multi-Head Self-Attention

<https://sebastianraschka.com/blog/2023/self-attention-from-scratch.html>

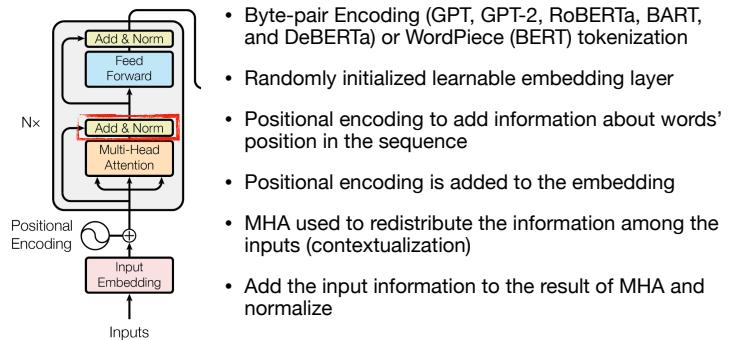


## Multi-Head Self-Attention

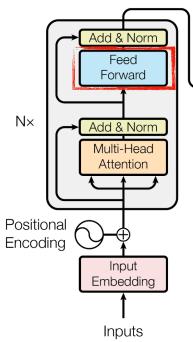
<https://sebastianraschka.com/blog/2023/self-attention-from-scratch.html>



## Transformer Encoder

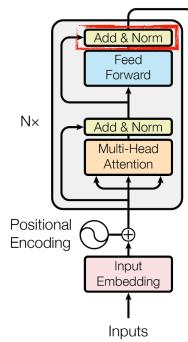


## Transformer Encoder



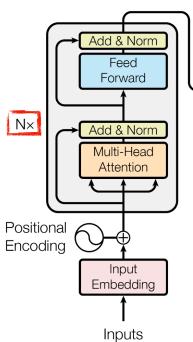
- Byte-pair Encoding (GPT, GPT-2, RoBERTa, BART, and DeBERTa) or WordPiece (BERT) tokenization
- Randomly initialized learnable embedding layer
- Positional encoding to add information about words' position in the sequence
- Positional encoding is added to the embedding
- MHA used to redistribute the information among the inputs (contextualization)
- Add the input information to the result of MHA and normalize
- Just a regular feed-forward network

## Transformer Encoder



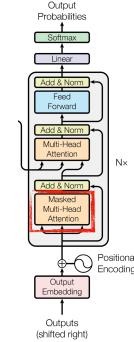
- Byte-pair Encoding (GPT, GPT-2, RoBERTa, BART, and DeBERTa) or WordPiece (BERT) tokenization
- Randomly initialized learnable embedding layer
- Positional encoding to add information about words' position in the sequence
- Positional encoding is added to the embedding
- MHA used to redistribute the information among the inputs (contextualization)
- Add the input information to the result of MHA and normalize
- Just a regular feed-forward network
- Another residual connection and layer normalization

## Transformer Encoder

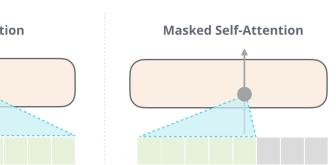


- Byte-pair Encoding (GPT, GPT-2, RoBERTa, BART, and DeBERTa) or WordPiece (BERT) tokenization
- Randomly initialized learnable embedding layer
- Positional encoding to add information about words' position in the sequence
- Positional encoding is added to the embedding
- MHA used to redistribute the information among the inputs (contextualization)
- Add the input information to the result of MHA and normalize
- Just a regular feed-forward network
- Another residual connection and layer normalization
- Repeat N times

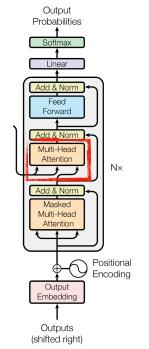
## Transformer Decoder



- Has the same elements as encoder
- Masked MHA is used so that the model cannot "look into the future"

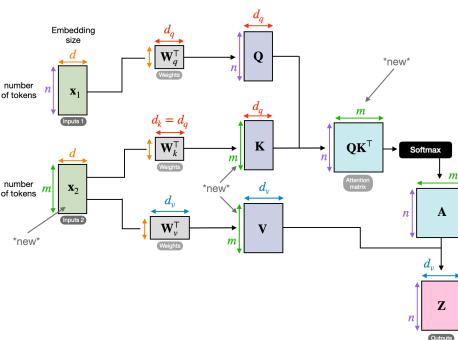


## Transformer Decoder

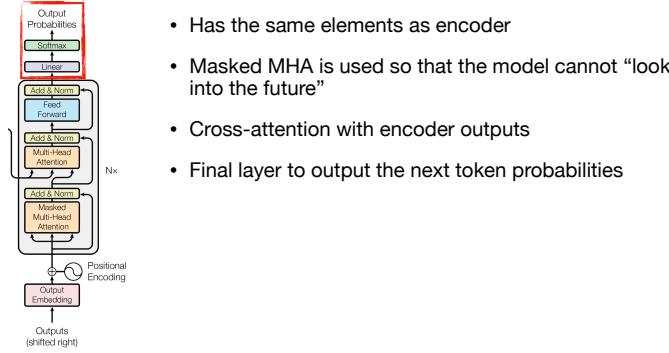


- Has the same elements as encoder
- Masked MHA is used so that the model cannot "look into the future"
- Cross-attention with encoder outputs

## Cross-Attention



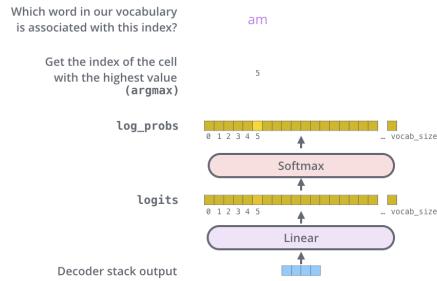
## Transformer Decoder



- Has the same elements as encoder
- Masked MHA is used so that the model cannot “look into the future”
- Cross-attention with encoder outputs
- Final layer to output the next token probabilities

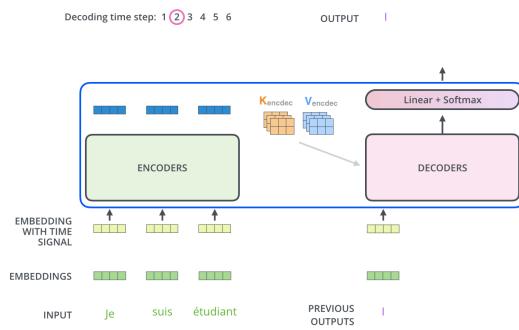
## Transformer Final Layer

<https://jalammar.github.io/illustrated-transformer/>



## Encoder-Decoder in Practice

<https://jalammar.github.io/illustrated-transformer/>

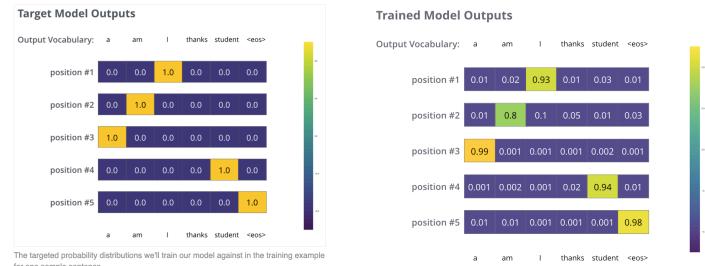


## Training

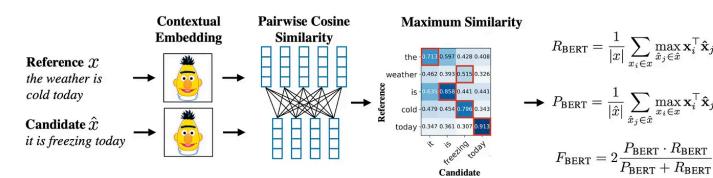
### Cross-Entropy Loss

$$H(P^* | P) = - \sum_i P^*(i) \log P(i)$$

TRUE CLASS DISTRIBUTION      PREDICTED CLASS DISTRIBUTION



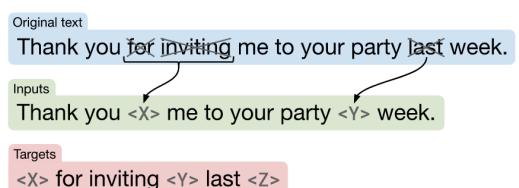
## BERT Score



Source: Bertscore: Evaluating text generation with bert

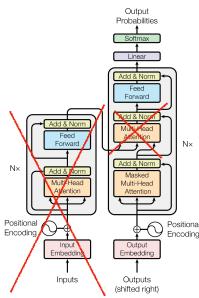
## T5

### Pre-training



## GPT-2

Radford, Alec, et al. "Language models are unsupervised multitask learners." (2019)



## Sentiment Analysis

Sentiment analysis (often referred to as opinion mining) is the process of gathering and analyzing people's opinions, thoughts and impressions regarding various topics, products, subjects and services.

- Rapid growth of internet based applications such as social media and blogs.
- People's opinions can be beneficial for
  - Corporations: product reviews on amazon, sentiment towards a certain feature or service.
  - Governments: public opinion on policies, expected voting tendencies before election.
  - Entertainment industry: movie reviews, game reviews.
  - Travel: reviews of restaurants, hotels, tourist attractions, etc.

## Sentiment Analysis

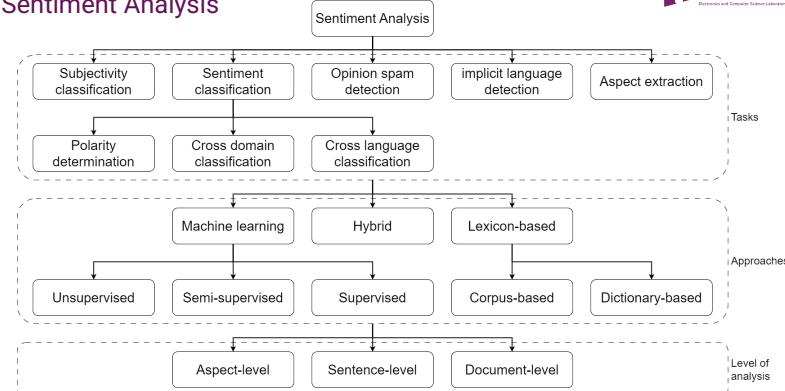


Sentiment analysis is a broad concept that consists of many different tasks, approaches and types of analysis.

Cambrai et al. argue that a holistic approach is required, and only classification or categorization is not sufficient. They present it as a 3 layered problem that includes 15 NLP problems:

- Syntactic layer: Microtext normalization, sentence boundary disambiguation, POS tagging, text chunking and lemmatization.
- Semantics layer: Word sense disambiguation, concept extraction, named entity recognition, anaphora resolution and subjectivity detection.
- Pragmatics layer: Personality recognition, sarcasm detection, metaphor understanding, aspect extraction and polarity detection.

## Sentiment Analysis



## Tasks



### □ Sentiment classification

- Most widely known and researched task in sentiment analysis.
- It can be divided into three major sub-tasks: polarity classification, cross-domain classification and cross-language classification.
- Polarity is usually classified as positive or negative with some researchers including a third category neutral.
- Cross-domain classification models transfer knowledge learned from data-rich source domain to target domain where data and/or labels are limited.
  - Extract domain invariant features whose distribution in the source domain is close to that in target domain
- Cross-language classification fulfills the same function but for languages.
  - An example can be to train the model in source language with abundant data and testing it on target language by translating the input to source language.

## Tasks

### □ Subjectivity classification

- The goal of subjectivity classification is to restrict unwanted objective data objects for further processing.
- It detects subjective clues, words that carry emotion or subjective notion like 'expensive', 'easy', 'better', etc.
- These clues are used to classify objects as subjective or objective.

Ce livre coûte 10 euros → Cannot be used for sentiment analysis

Ce livre est cher → Can be used for sentiment analysis

## Tasks

### Opinion spam detection

- Opinion spams refer to fake or false reviews intelligently written to either promote or discredit a product.
- Three main features are considered within this context:
  - Review content: the actual text of the review
  - Meta-data: information like IP address, geo-location, user-id, etc.
  - Real-life knowledge: this method utilizes learned experiences to classify spam. For example, if a product has good reputation and suddenly inferior ratings are given over a period, reviews of that period might be suspected.

## Tasks

### Implicit language detection

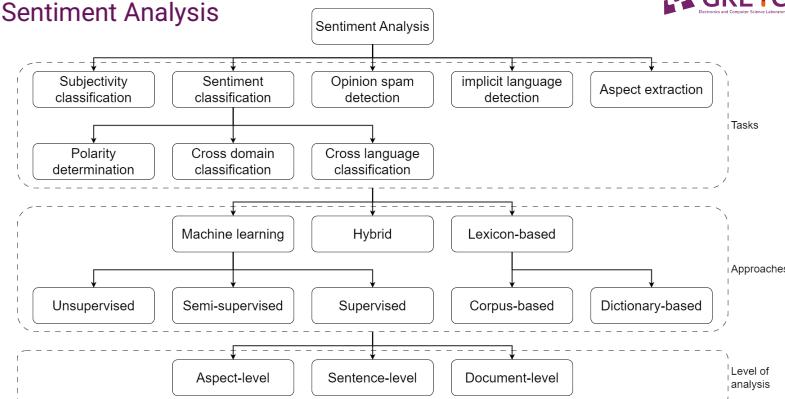
- Implicit language includes humor, sarcasm and irony.
- There is vagueness and ambiguity in this form of speech, which is sometimes hard to detect.
- An implicit meaning can sometimes completely flip the polarity of a sentence.
- Example, « I love pain », pain is a factual word with negative polarity. The contradiction between pain and love can indicate sarcasm.
- Traditional methods for detection include exploring emoticons, expressions of laughter and heavy punctuation mark usage.

## Tasks

### Aspect extraction

- Aspect extraction refers to retrieving the target entity and the aspects of the target entity in the document. The target entity can be a person, product, event, etc.
- Aspect extraction is particularly important in sentiment analysis of social media and blogs that often do not have predefined topics.
- Most traditional method for this is frequency based where most frequent nouns and compound nouns are considered as candidates for aspects.
  - Not all nouns are aspects
- Syntax-based methods find aspects by means of syntactic relations they are in. For example, identifying aspects that are preceded by a modifying adjective that is a sentiment word.
  - Many relations need to be found for complete coverage

## Sentiment Analysis



## Lexicon based approaches

- Traditional approach for sentiment analysis that scans through the text for words that express positive or negative feelings to humans.
- It shows to be extremely dependent on domain of interest due to differences in language usage between domains.
- There are two main approaches to creating sentiment lexicons: dictionary based and corpus based.
- Dictionary based approaches start with an initial list of terms and iteratively expand the lexicon by adding synonyms and antonyms of the terms to the list.
- They work best for general purpose use.
- Corpus based approaches starts with general purpose list of words and finds other terms from domain specific corpus based on co-occurring word patterns.

## Image Captioning



## Global CNN features

With advent of CNNs, all models consuming visual inputs have been improved in terms of performance

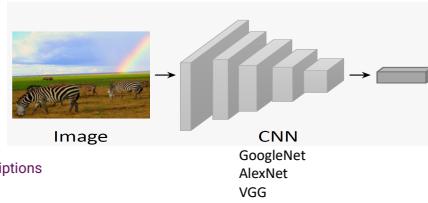
In the most simple recipe, the activation of one of the last layers of a CNN is employed to extract high-level representations, which are then used within language models for generating final text.

Advantages

- Simplicity
- Compactness of representations

Disadvantages

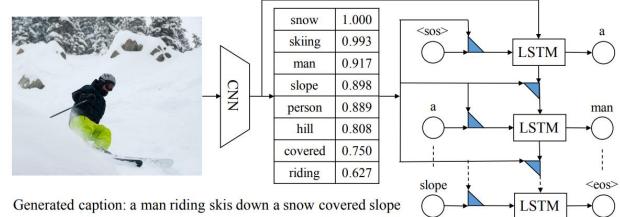
- Excessive compression of information
- Lack of granularity
- Unable to produce specific and fine-grained descriptions



## Global CNN features

Gan et al., CVPR 2017

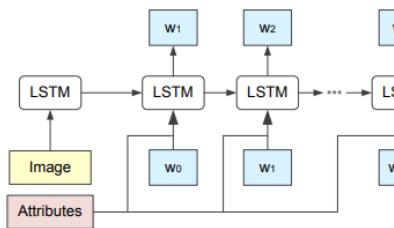
- ResNet-152 pre-trained on imagenet dataset



## Global CNN features

Yao et al., Boosting image captioning with attributes, ICCV 2017

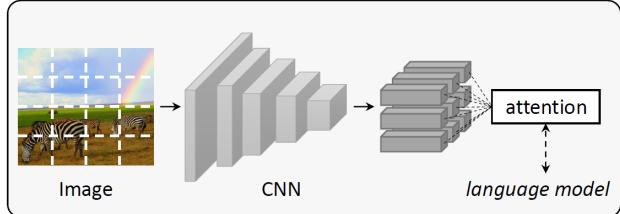
- 1,000 most common words on COCO as the high-level attributes and train the attribute detectors with MIL model (Fang et al., 2015)
- GoogleNet for image encoding



## Attention based models

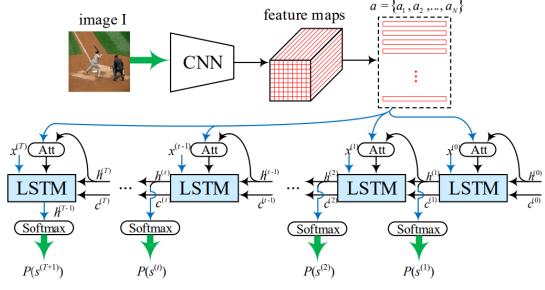
Grid based attention

- Improves upon the drawbacks of global representations and increases the granularity level of visual encoding
- Motivated from use of attention in machine translation.



## Attention based models

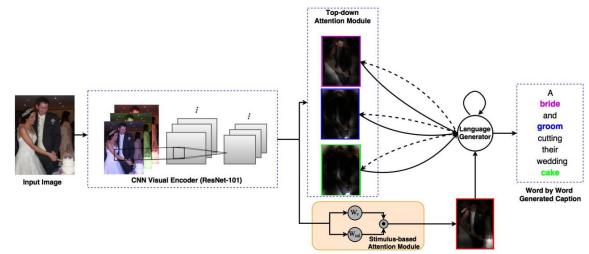
Ge et al., Exploring Overall Contextual Information for Image Captioning in Human-Like Cognitive Style, ICCV 2019



## Attention based models

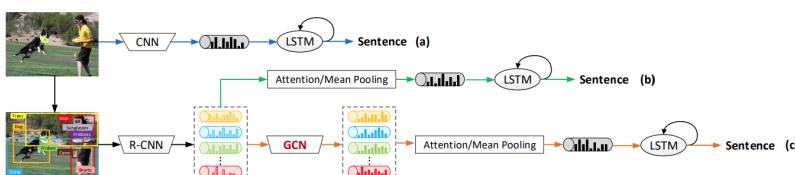
Chen et al., Boosted attention: leveraging human attention for image captioning, ECCV, 2018

- Improve the attention mechanism by incorporating human attention input into the model.



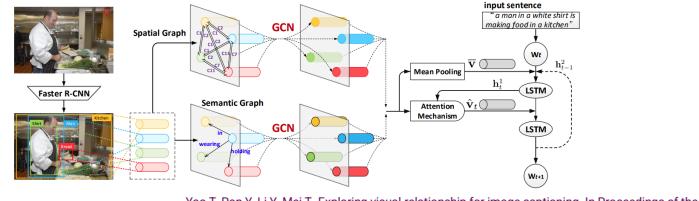
## Graph based models

- Region based attention treats all regions equally without taking into account the interactions between them.
- Some studies consider using graphs based models for improved encoding of image regions by incorporating relations between different regions.



## Graph based models

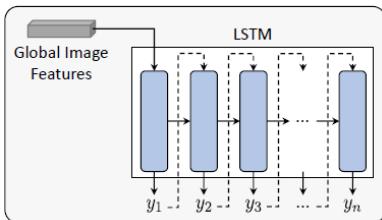
- Spatial and semantic graphs
- A semantic relation classifier is trained on a large corpus and directly used for semantic graph generation.
- Spatial graphs are built and assigned depending on their Intersection over Union (IoU), relative distance and angle.



Yao T, Pan Y, Li Y, Mei T. Exploring visual relationship for image captioning. In Proceedings of the European conference on computer vision (ECCV) 2018 (pp. 684-699).

## Language models

- The main language modeling strategies applied to image captioning are:
  1. LSTM based
  2. Transformer based fully attentive approaches
  3. Image-text early fusion (BERT-like)



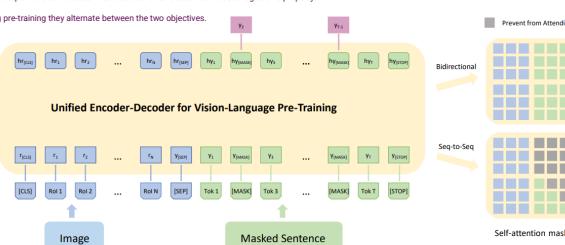
## Vision-Language Transformer Model

- Unified the transformer encoder and decoder into a single model.
- Model input consists of class-aware region embeddings, word embeddings and three special tokens [CLS], [SEP] and [STOP].
- [CLS] indicates the start of the visual input.
- [SEP] indicates the separation between visual and text input.
- [STOP] indicates end of sentence.
- [MASK] indicates the masked word



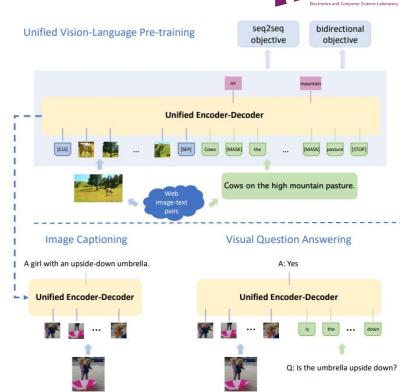
## Vision-Language Transformer Model (pre-training)

- 15% of the text tokens are replaced by [MASK] token, random token or original token.
- The hidden state from the last transformer block is projected to word likelihoods where the masked token is predicted as classification problem.
- Through this reconstruction the model learns the dependencies in the context and forms a language model.
- Two objectives are considered within this model:
  - Bidirectional: every token can attend to every other token.
  - Seq-to-seq: tokens cannot attend to future tokens. It satisfies the auto-regressive property.
- During pre-training they alternate between the two objectives.



## VLP model for image-captioning

- For image captioning we fine-tune using seq2seq objective.
- During inference
  - Encode image region along with special tokens ([CLS] and [SEP] tokens).
  - We then start the generation by feeding in the [MASK] token and sampling a word from word likelihood output.
  - Replace the [MASK] token in the input sequence with sampled word and add new [MASK] token to the end of sequence.
  - Generation terminates when [STOP] token is chosen.



## Training strategies

### Cross-Entropy Loss

- Most used objective for image captioning.
- The aim is to minimize the negative log-likelihood of the current word given the previous ground-truth words.
- The loss works at word level and optimizes the probability of each word without considering long range dependencies.

$$L_{XE}(\theta) = - \sum_{i=1}^n \log (P(y_i | y_{1:i-1}, \mathbf{X}))$$

## Training strategies

### Masked Language Model

- Idea is to randomly mask a subset of input tokens and train the model to predict these masked tokens based on remaining tokens, both previous and subsequent.
- The model relies more on context making it more robust.
- Training on these models is much slower since they only train in masked tokens and not entire sentence.

## Evaluation

### BLEU score

- Target sentence: The guard arrived late because it was raining
  - Predicted sentence: The guard arrived late because of the rain
- We first calculate precision scores for 1-gram through 4-grams.

**Target Sentence:** The guard arrived late because it was raining  
**Predicted Sentence:** The guard arrived late because of the rain  
 $P_1 = 5/8$

**Target Sentence:** The guard arrived late because it was raining  
**Predicted Sentence:** The guard arrived late because of the rain  
 $P_2 = 4/7$

**Target Sentence:** The guard arrived late because it was raining  
**Predicted Sentence:** The guard arrived late because of the rain  
 $P_3 = 3/6$

**Target Sentence:** The guard arrived late because it was raining  
**Predicted Sentence:** The guard arrived late because of the rain  
 $P_4 = 2/5$

## Evaluation

### BLEU score

- Brevity penalty: it penalizes sentences that are too short.

If the predicted sentence is just « the », the 1-gram precision is 1/1=1, indicating perfect score.

$$\text{Brevity Penalty} = \begin{cases} 1, & \text{if } c > r \\ e^{(1-r/c)}, & \text{if } c \leq r \end{cases}$$

c = predicted sentence length  
r = target sentence length

$$\text{Bleu}(N) = \text{Brevity Penalty} \cdot \text{Geometric Average Precision Scores}(N)$$

$$\text{Geometric Average Precision}(N) = (p_1)^{\frac{1}{N}} \cdot (p_2)^{\frac{1}{N}} \cdot (p_3)^{\frac{1}{N}} \cdot (p_4)^{\frac{1}{N}}$$

## Evaluation

### ROUGE

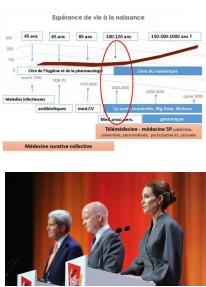
- Compared to BLEU score that focuses on precision, ROUGE focuses on recall.

### ROUGE-N

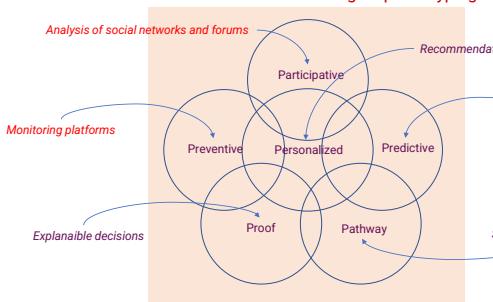
$$\text{ROUGE-N} = \frac{\sum_{S \in \{\text{ReferenceSummaries}\}} \sum_{\text{gram}_n \in S} \text{Count}_{\text{match}}(\text{gram}_n)}{\sum_{S \in \{\text{ReferenceSummaries}\}} \sum_{\text{gram}_n \in S} \text{Count}(\text{gram}_n)}$$

## 6P Medicine

- 1P - Personalized:** Personalized medicine consists of adapting a medical treatment according to the individual characteristics of a patient.
- 2P - Preventive:** Preventive medicine focuses on wellness, and consists of measures taken for disease prevention.
- 3P - Predictive:** Predictive medicine is a branch of medicine that aims to identify patients at risk of developing a disease.
- 4P - Participative:** Medicine should be participatory, leading patients to be more responsible for their health and care.
- 5P - Proof:** Medicine must be based on evidence of medical service to patients, especially when it relies on connected health and telemedicine.
- 6P - Pathway:** Coordinating multiple interventions (medical, social, occupational medicine, etc.) such that the healthcare pathway is progressively articulated, according to the pathology and its evolution.

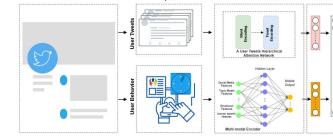


1



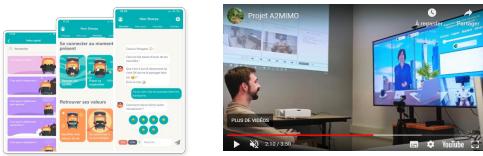
## Social Network Analysis

- Social networks are an important support for Participative medicine, which automatic analysis might allow Preventive/Predictive actions.
- It is common for people who suffer from mental health problems to often disclose their feelings and their daily struggles with mental health issues on social media as a way of relief.
- Twitter, Reddit, Doctissimo, to name but a few platforms have become an excellent resource to automatically discover people who are under depression.
- [Zogan et al., 2021] propose a depression detection framework by tackling textual, behavioral, temporal, and semantic modalities.



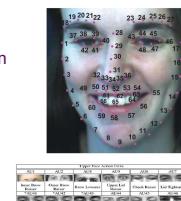
## Monitoring Platforms and (Embodied) Chatbots

- A chatbot is a system that is able to converse and interact with human users using spoken, written, and visual languages (embodied).
- Chatbots can be useful preventive tools for individuals who are reluctant to seek mental health advice due to stigmatization.
- [Abd-alrazaq et al., 2019] studied 41 different embodied and non-embodied chatbots. Most tackle depression and autism.
- Among other scientific issues, therapeutic alliance is the key factor for the success of chatbots and ECAs.



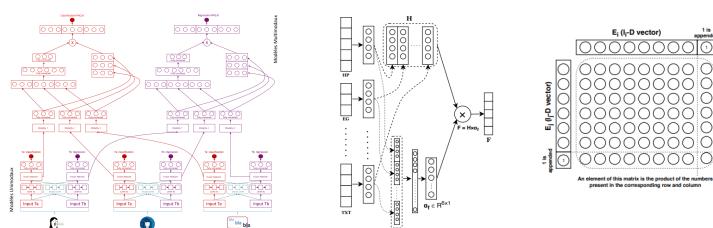
## Multimodal Estimation of PHQ-8

- In a patient-therapist interview, different signals should be combined for a correct diagnosis.
- Within the DAIC-WOZ dataset, the following signals are available:
  - Visual signals : expression of sadness, gaze escape, etc.
    - Facial Landmarks (FL), Head Pose (HP), Eye Gaze (EG), Action Unit (AU).
  - Speech signals : veiled voice, monotonous tone, etc.
    - Formant (FMT), COVAREP (COV).
  - Language signals : negative vocabulary, lack of perspective, etc.
    - Universal Sentence Encoder (TR).



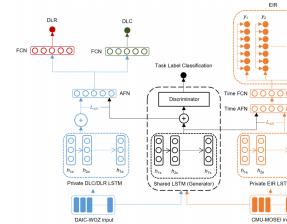
## Multimodal Estimation of PHQ-8

- Combining classification and regression of depression estimators.
- An attention fusion network is used to combine inputs.
- Intra-modality inputs signals are combined with tensors.



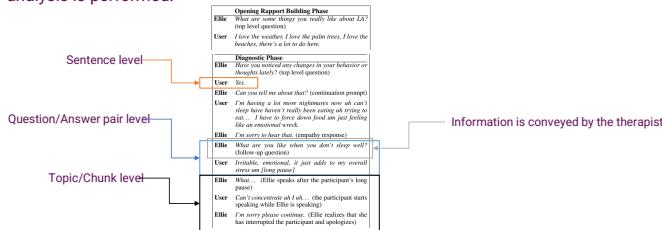
## Emotional Language for the Estimation of PHQ-8

- In [Qureshi et al., 2020], we hypothesize that the estimation of depression level can benefit from the concurrent learning of emotion intensity.
- The CMU-MOSEI dataset comprises 3,228 videos from 1,000 different speakers over 250 topics. Videos were gathered from an online video platform, where users emit their opinions in the form of monologues.



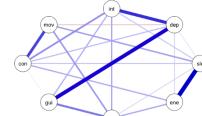
## Analysis of Structured Interviews

- First observation: most of the related works have been dealing with the interview on a **line by line basis**; the hypothesis being that sentence representation is the correct one.
- Second observation: some of the related works only deal with the patient information; the hypothesis being that **only the patient information is important** for the diagnosis.
- Our hypothesis is that better diagnosis can be established if the correct level of language analysis is performed.



## Symptom-based Analysis

- Most related works have been tackling depression level estimation as a **simple task** (depressed or non depressed). More advanced models have been trying to predict the **PHQ-8 score** (between 0 and 24) directly or propose to solve the **intermediate 5-class problem** (none-minimal, mild, moderate, moderately severe, severe depression).
- In Psychiatry, there is a **shift towards richer representations of psychiatric syndromes** that can take into account the dimensional and heterogeneous nature of the clinical pictures of the same psychiatric diagnosis. One particular approach that is gaining attention concerns **symptom network analysis**.
- We develop similar models as previously to acknowledge if they can handle the **prediction of individual symptom values**, where each of the 8 symptoms is a value between 0 and 3.



## Symptom-based Analysis

- In order to better understand results, we present a **radar plot analysis** that shows that adequate behavior of the model is obtained.

