

# Επεξεργασία Φυσικής Γλώσσας

# Distributed Contextual Embeddings

Demetris Paschalides

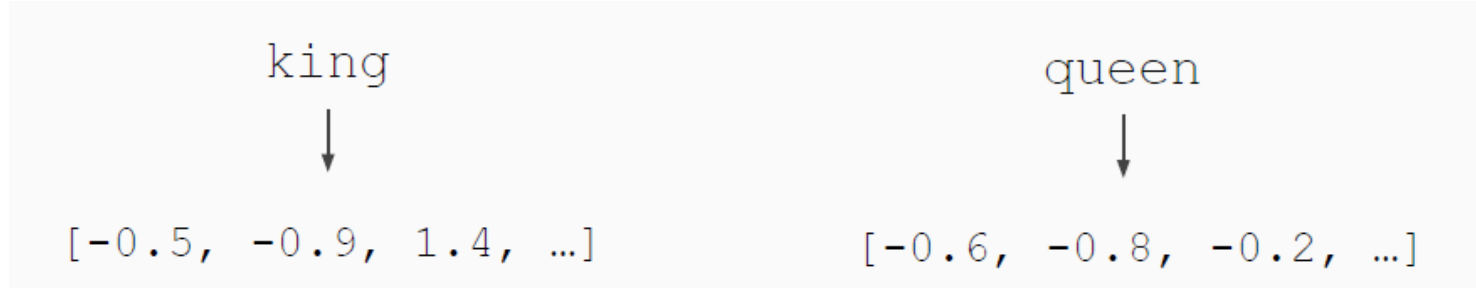
Department of Computer Science

University of Cyprus



# Word Embeddings

- Word embeddings αποτελούν τη βάση της βαθιάς μάθησης για NLP



- Word embeddings (*word2vec*, *GloVe*) είναι συχνά προ-εκπαιδευμένοι σε σώματα κειμένων από στατιστικά στοιχεία





# Παραστάσεις Word Embedding

- ❑ Count-based
  - TF-IDF, PPMI
- ❑ Class-based
  - Brown Clusters
- ❑ Distributed Prediction-based Embeddings
  - Word2Vec, FastText
- ❑ Distributed Contextual Embeddings from Language Models
  - Elmo, BERT
- ❑ Variants
  - Multi-lingual, Multi-sense etc.



# Context είναι τα πάντα

- $p(\textit{play} \mid \textit{Elmo and Cookie Monster play a game.})$
- $p(\textit{play} \mid \textit{The Broadway play premiered yesterday.})$

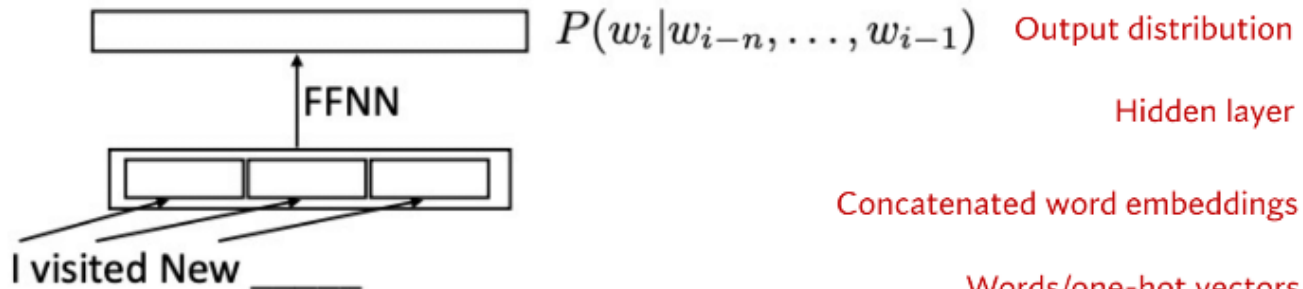


- Πώς να χειριστείτε διαφορετικές έννοιες λέξεων?
  - Μέχρι στιγμής, τα μοντέλα ενσωμάτωσης λέξεων παράγουν ένα διάνυσμα για τη λέξη “play”.
- Εκπαιδεύστε ένα **neural language model** για να προβλέψετε την επόμενη λέξη που δίνεται προηγούμενες λέξεις στην πρόταση

# Neural Language Models

- ❑ Θυμηθείτε το Language Modeling task.
- ❑ **Input:** ακολουθία λέξεων.
- ❑ **Output:** πιθανότητα της επόμενης λέξης  $w$ .
  
- ❑ Αρχικές προσεγγίσεις: Feed Forward Neural Networks (FFNN) κοιτάζοντας context.

$$P_{Add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$



# Neural Language Models

- ❑ Βελτιώσεις σε σχέση με N-Gram Language Model:
  - Δεν υπάρχει πρόβλημα αραιότητας.
  - Δεν χρειάζεται να αποθηκεύσετε όλα τα παρατηρούμενα N-Grams.
- ❑ Περιορισμούς:
  - Το σταθερό παράθυρο είναι πολύ μικρό.
  - Μεγέθυνση παραθύρου μεγεθύνει  $W$ .
  - Τα Windows δεν μπορούν ποτέ να είναι αρκετά μεγάλα!
  - Διαφορετικές λέξεις πολλαπλασιάζονται με εντελώς διαφορετικά βάρη. Δεν υπάρχει συμμετρία στον τρόπο επεξεργασίας των εισόδων.



# Neural Language Models

- ❑ Improvements over N-Gram Language Model:
  - No sparsity problem.
  - Don't need to store all observed N-Grams.

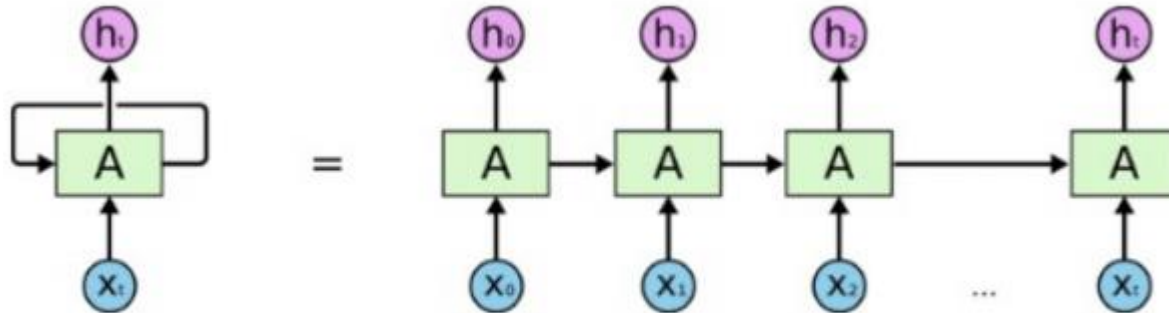
- ❑ Limitations:
  - Fixed window is too small.
  - Enlarging window enlarges  $W$ .
  - Windows can never be large enough!
  - Different words are multiplied by completely different weights. No symmetry in how the inputs are processed.

We need a neural architecture that can process any length.





# Recurrent Neural Networks (RNN)

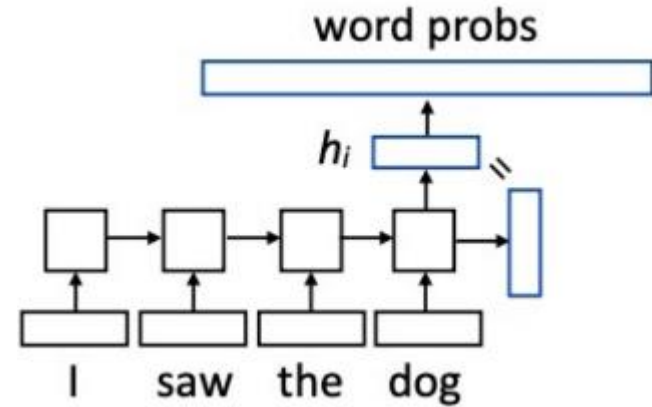


An unrolled recurrent neural network.

- Πάρτε διαδοχική είσοδο οποιουδήποτε μήκους.
- Εφαρμόστε τα ίδια βάρη σε κάθε βήμα.
- Προαιρετικά παράγετε έξοδο σε κάθε βήμα.

# RNN Language Modeling

- $P(w \mid \text{context}) = \text{softmax}(Wh_i)$
- $W$  είναι ένα matrix (μέγεθος λεξιλογίου)  $\times$  (hidden size).

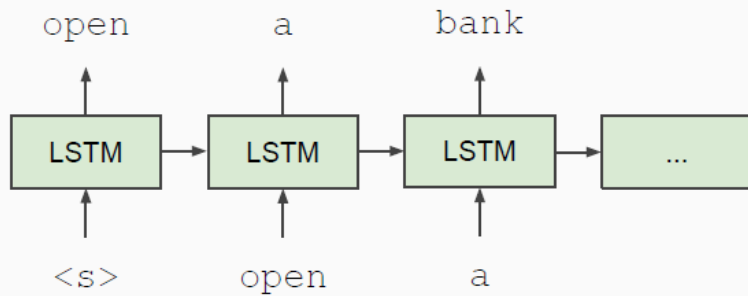


## Training RNN Language Model

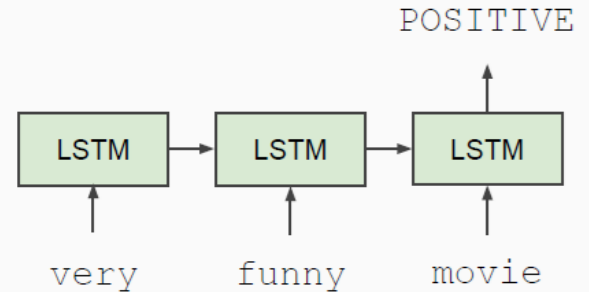
- Input = μια ακολουθία λέξεων
- Output = Οι λέξεις άλλαξαν μία θέση
- Αυτό μας επιτρέπει να ομαδοποιούμε αποτελεσματικά την εκπαίδευση με την πάροδο του χρόνου.
- Επίσης, παράγει χαρακτηριστικά από όλες τις προτάσεις / παραγράφους.

# Ημι-εποπτευόμενη μάθηση ακολουθίας

## Train LSTM Language Model

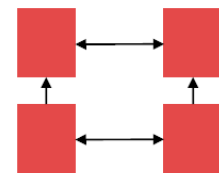


## Fine-tune on Classification Task



# Ζήτημα Linear Interaction Distance

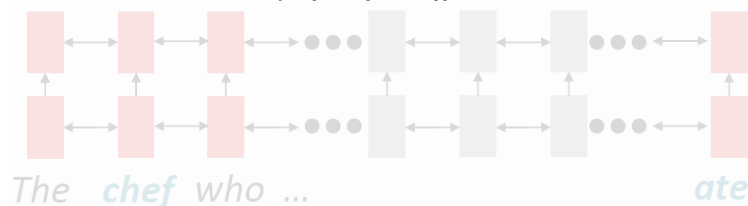
- ❑ RNNs ξετυλίγονται "από αριστερά προς τα δεξιά".
- Κωδικοποιεί γραμμική τοπικότητα: μια χρήσιμη ευρετική;



*tasty pizza*

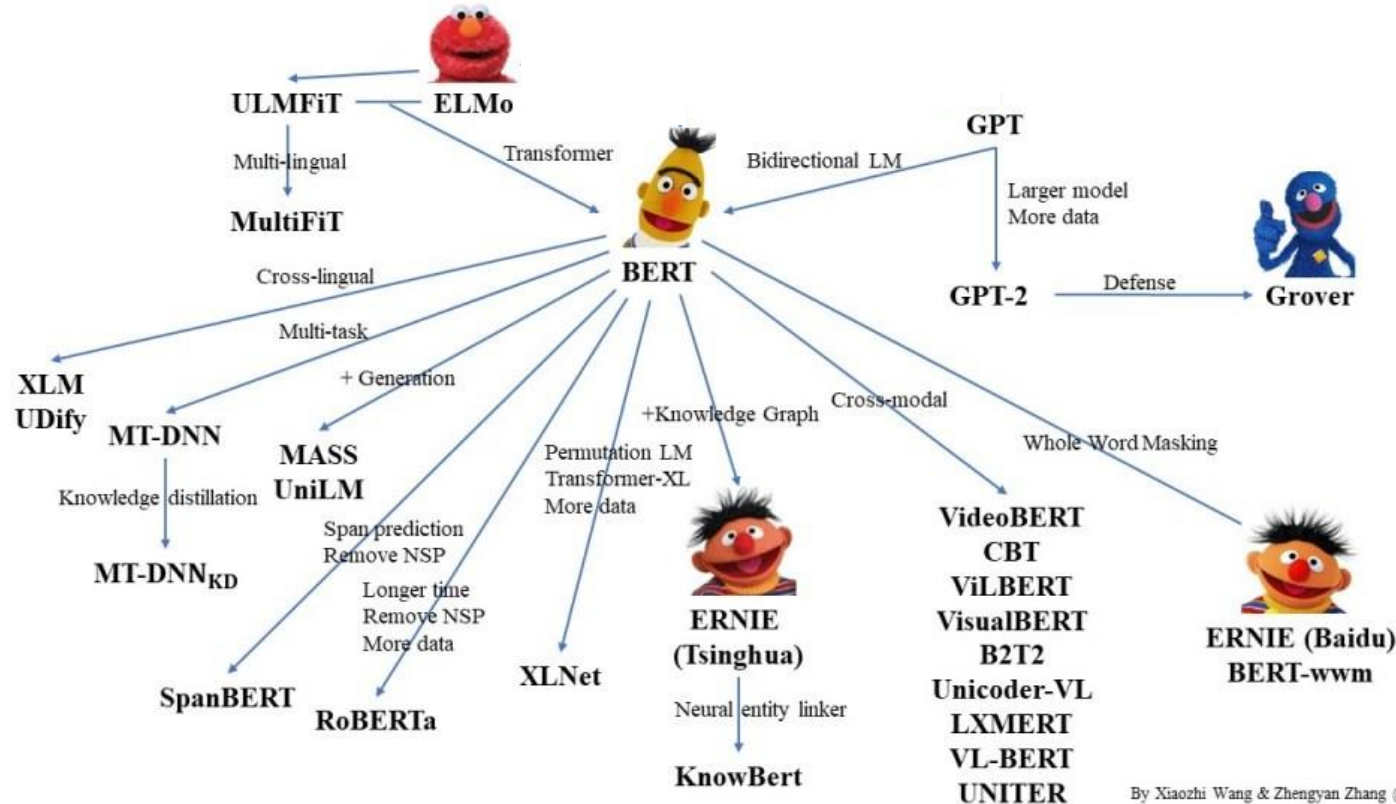
## Linear Interaction Distance:

- ❑  $O(\text{sequence length})$  Βήματα για την αλληλεπίδραση μακρινών ζευγών λέξεων σημαίνει:
  - Δύσκολο να μάθουν εξαρτήσεις μεγάλων αποστάσεων (προβλήματα κλίσης);
  - Η γραμμική σειρά των λέξεων "ψήνεται";
- ❑ **Περιορισμούς:** Χρειάζεστε μηχανισμό κατάδειξης για να επαναλάβετε πρόσφατες λέξεις.



- ❑ **Solution:** Transformers

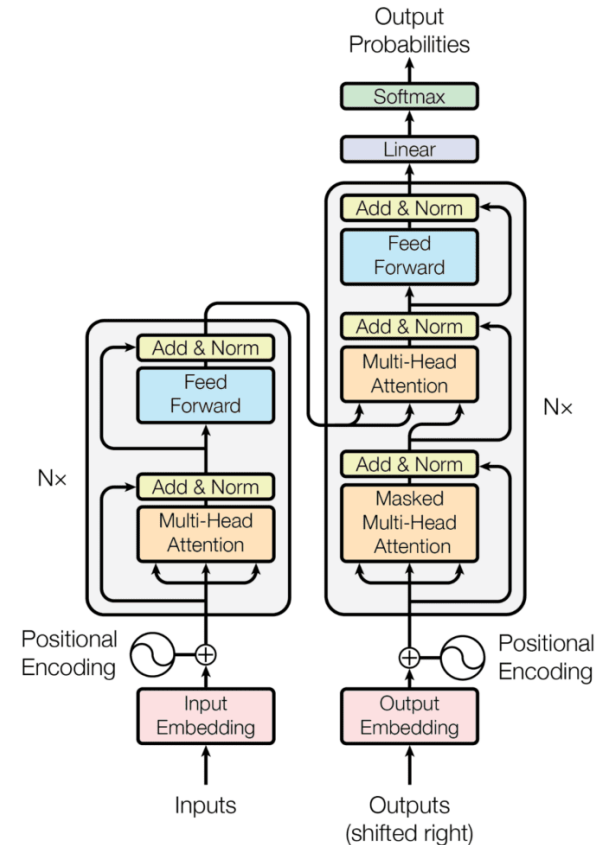
# Ιστορία στις Transformers



By Xiaozhi Wang & Zhengyan Zhang @THUNLP

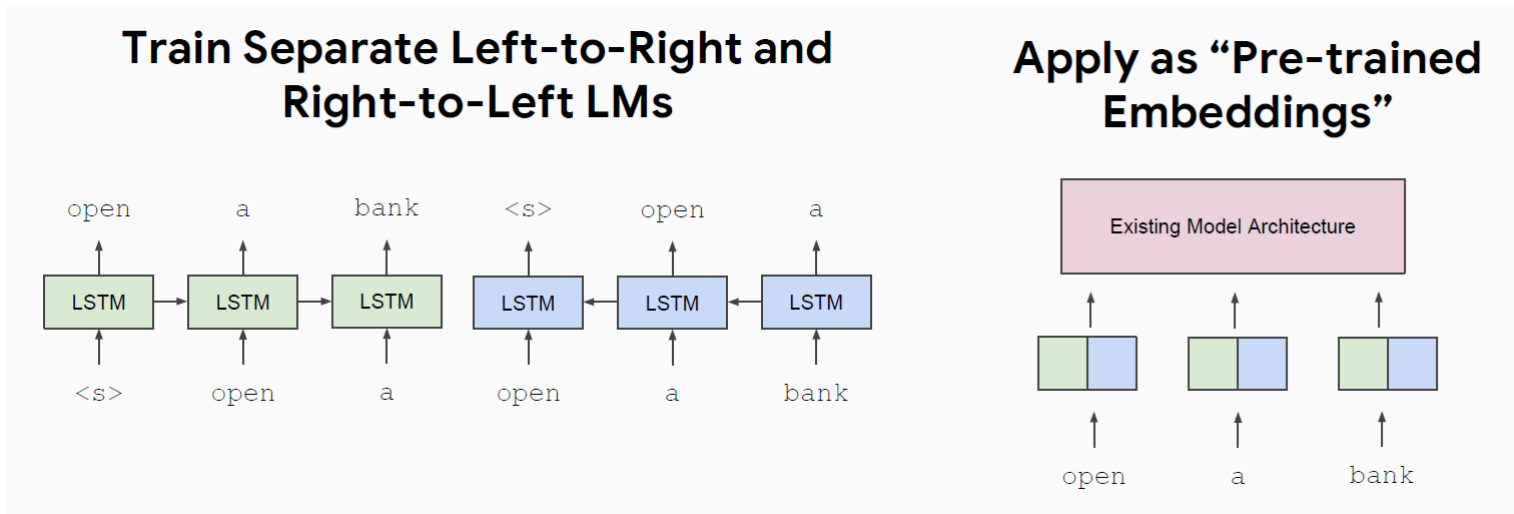
# Αρχιτεκτονική Transformer

- ❑ Η αρχιτεκτονική μετασχηματιστών ακολουθεί μια δομή κωδικοποιητή-αποκωδικοποιητή.
- ❑ Ο κωδικοποιητής αντιστοιχίζει μια ακολουθία εισόδου σε μια ακολουθία συνεχών αναπαραστάσεων, η οποία στη συνέχεια τροφοδοτείται σε έναν αποκωδικοποιητή.
- ❑ Ο αποκωδικοποιητής λαμβάνει την έξοδο του κωδικοποιητή μαζί με την έξοδο του αποκωδικοποιητή στο προηγούμενο χρονικό βήμα για να δημιουργήσει μια ακολουθία εξόδου.



# ELMo: Deep Contextualized Word Representation

- Instead of using a fixed embedding for each word, ELMo looks at the entire sentence before assigning each word in it an embedding.



# A Review of ELMo

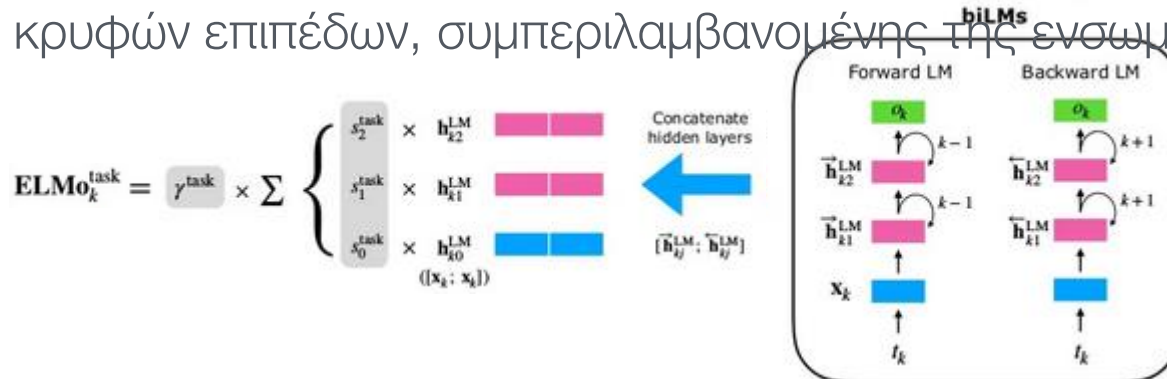
□ Το ELMo μαθαίνει word embeddings μέσω της δημιουργίας αμφίδρομων LM (BiLMs).

- BiLMs consist of forward and backward LMs.

- Forward:  $p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k | t_1, t_2, \dots, t_{k-1})$

- Backward:  $p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k | t_{k+1}, t_{k+2}, \dots, t_N)$

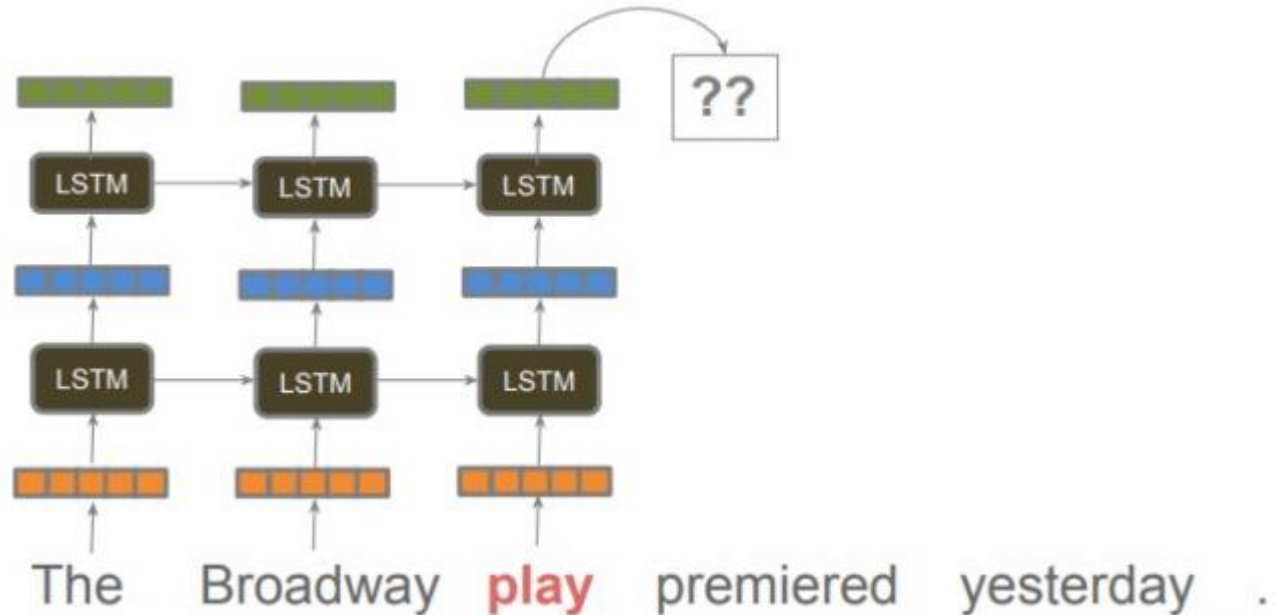
□ ELMo αντιπροσωπεύει μια λέξη  $t_k$  ως γραμμικός συνδυασμός αντίστοιχων κρυφών επιπέδων, συμπεριλαμβανομένης της ενσωμάτωσής του.





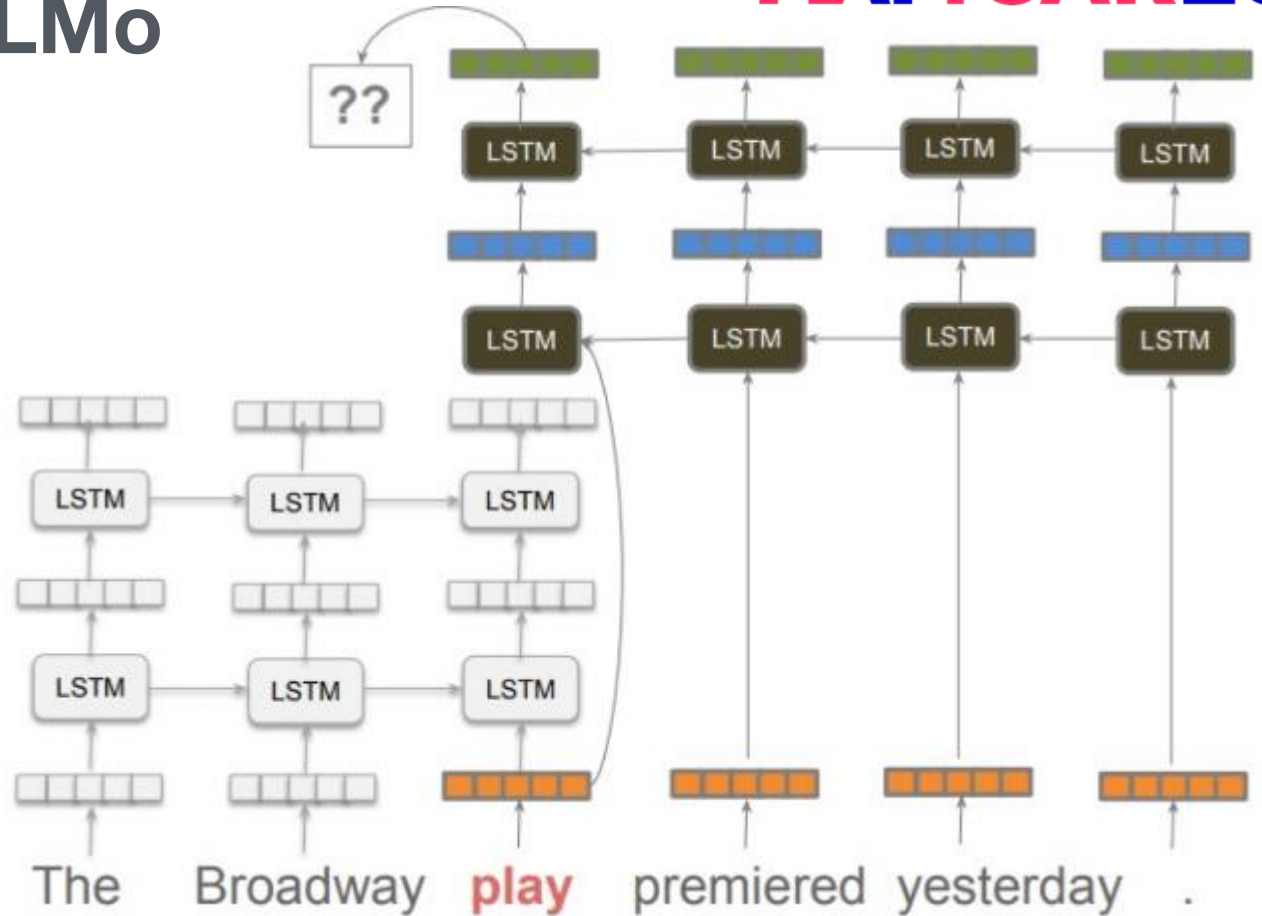
# A Review of ELMo

- An example of ELMo:



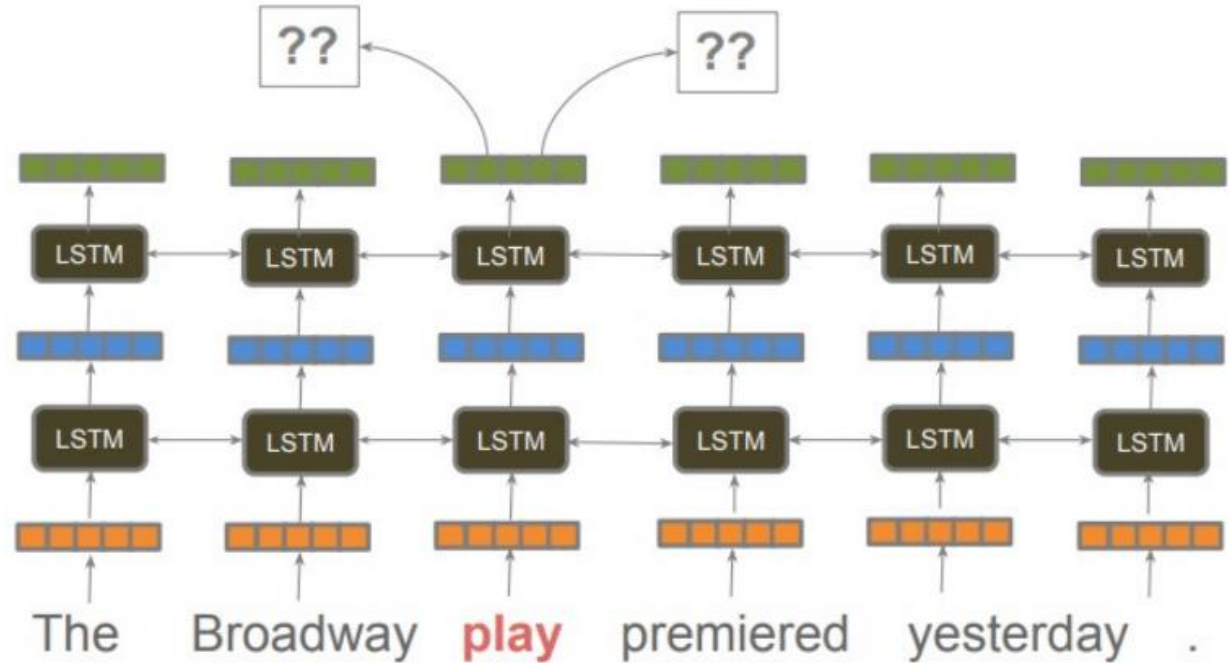
## A Review of ELMo

- An example of ELMo:



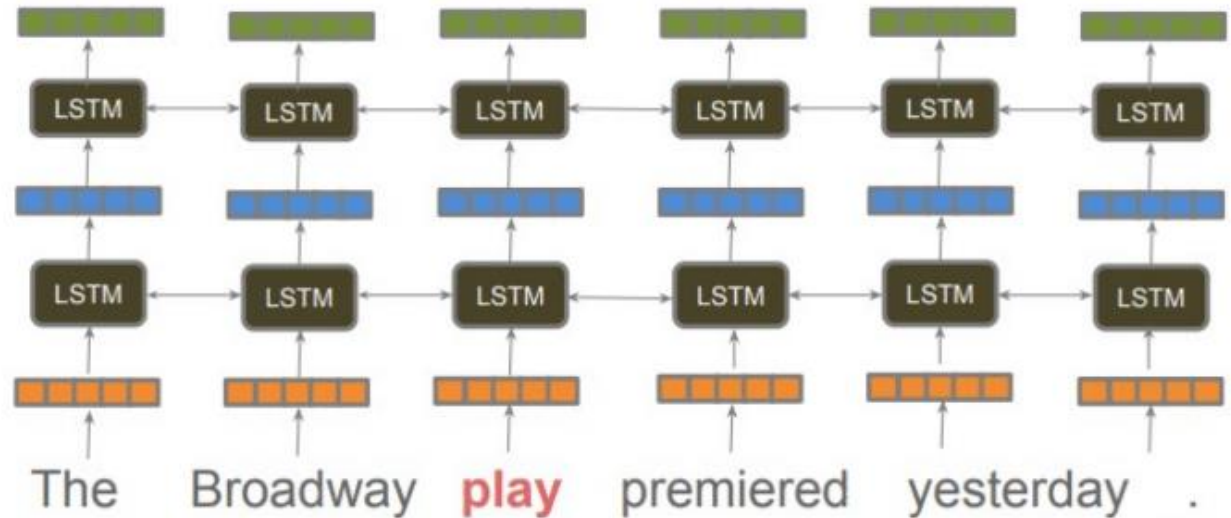
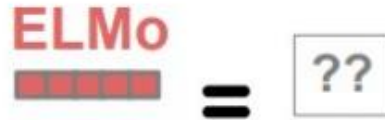
# A Review of ELMo

- An example of ELMo:



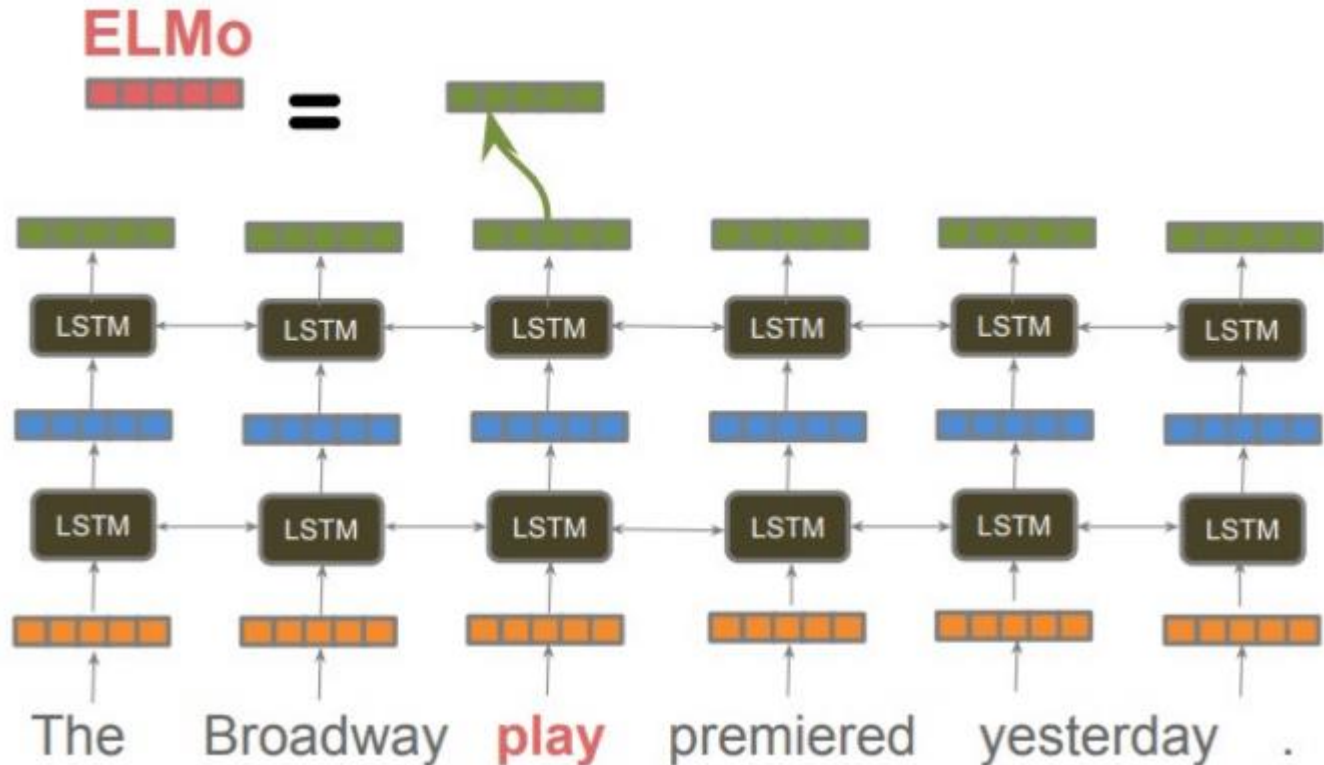
# A Review of ELMo

- An example of ELMo:



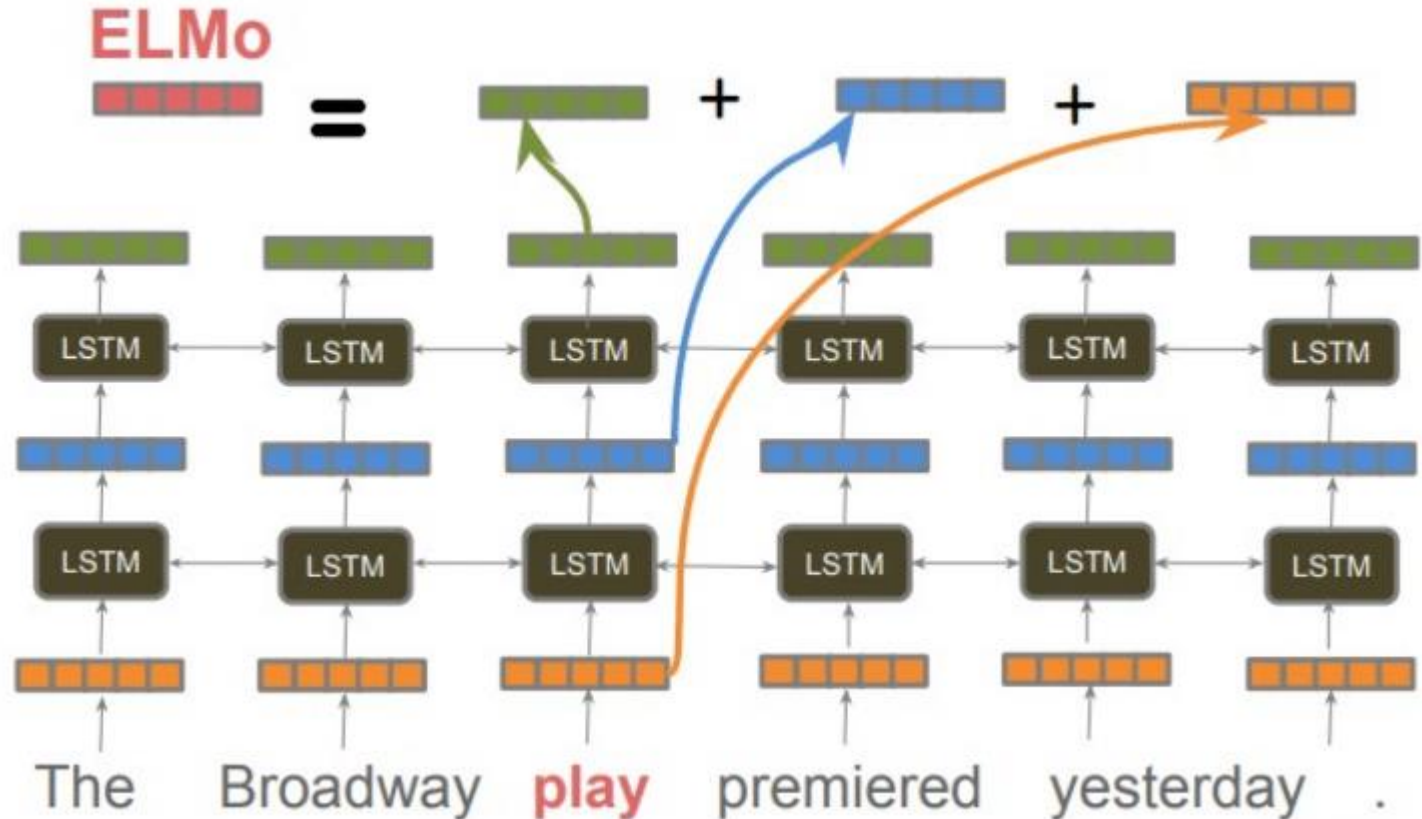
# A Review of ELMo

- An example of ELMo:



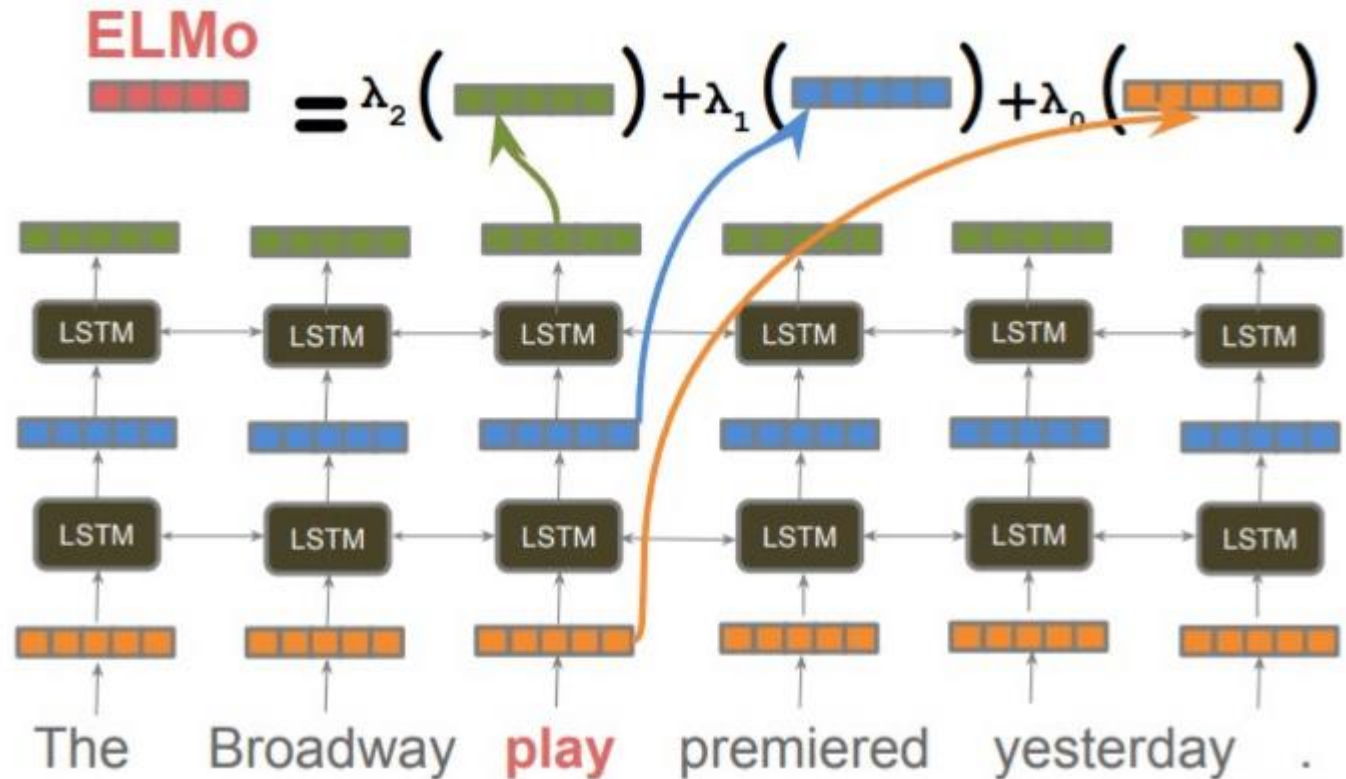
# A Review of ELMo

- An example of ELMo:



# A Review of ELMo

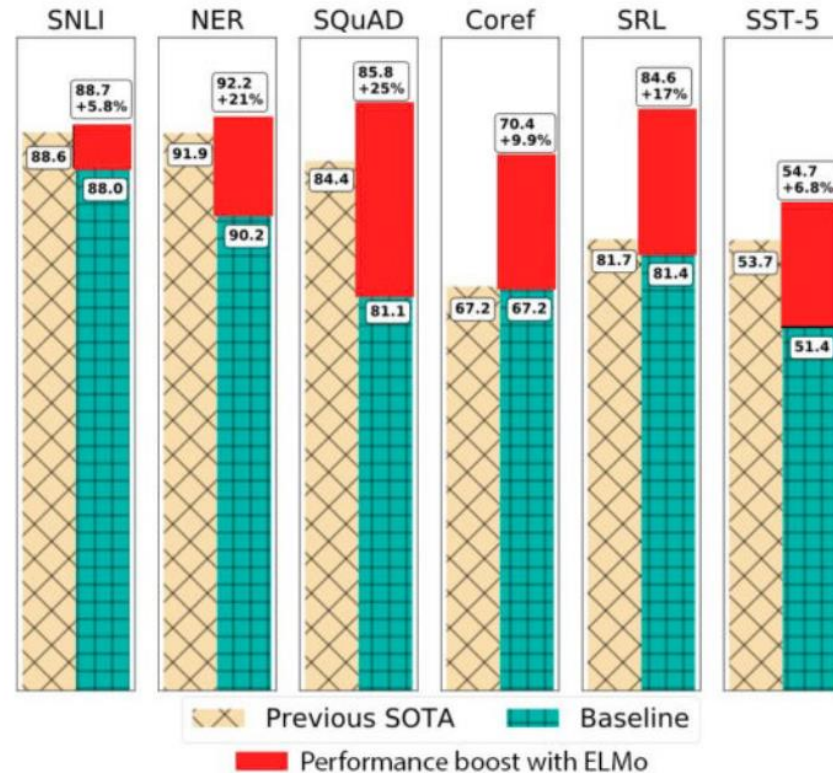
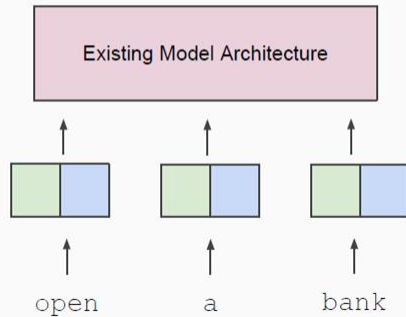
- An example of ELMo:





# Αξιολόγηση με εξωγενή καθήκοντα

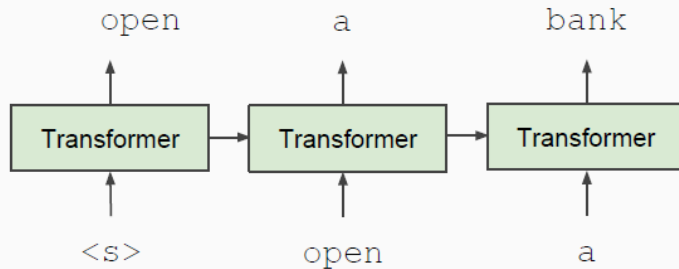
Apply as “Pre-trained Embeddings”



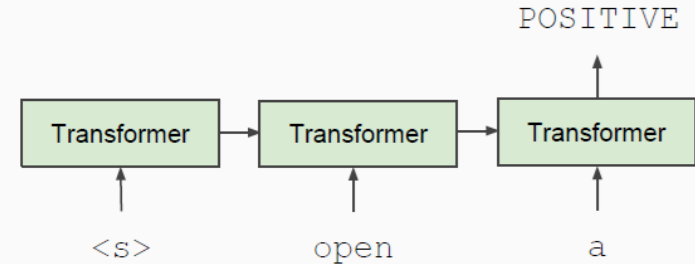


# OpenAI GPT: Generative Pre-trained Transformer

## Train Deep (12-layer) Transformer LM



## Fine-tune on Classification Task



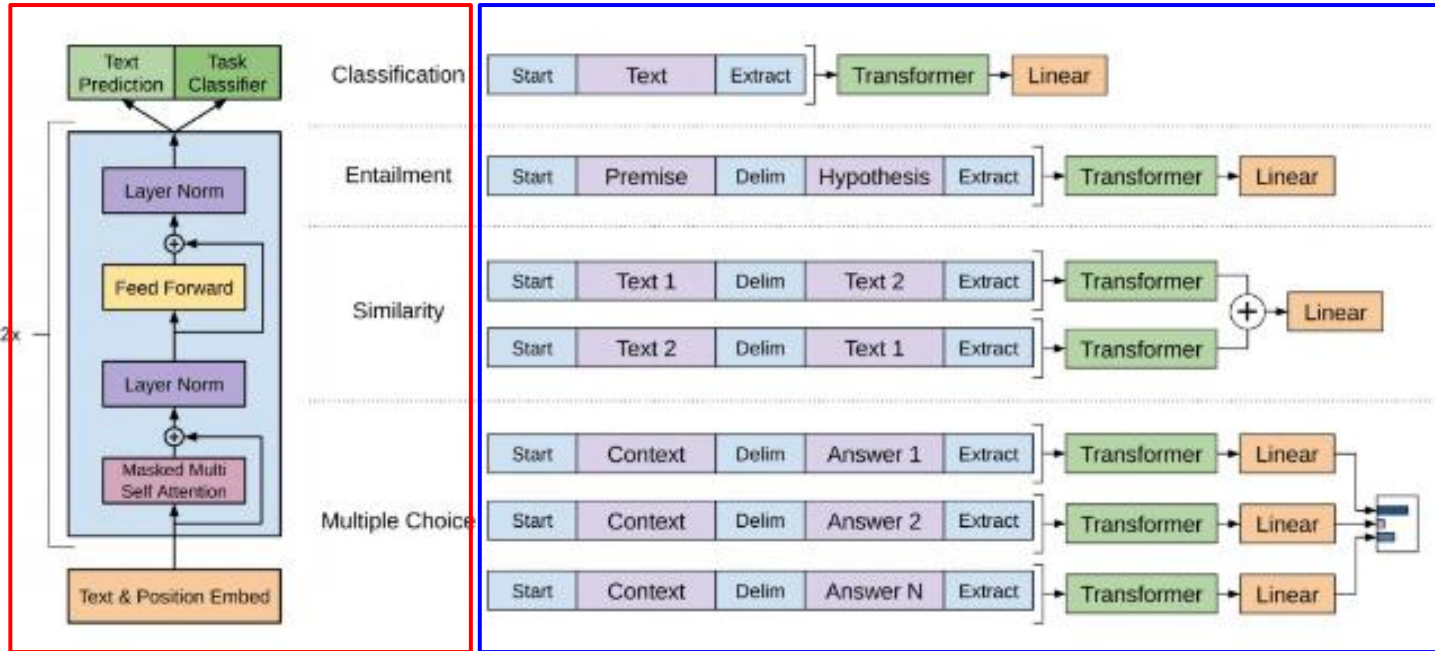
# OpenAI GPT: Generative Pre-trained Transformer

- Προ-εκπαίδευση χωρίς επίβλεψη για μεγιστοποίηση της log-likelihood :

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$

Πού  $\mathbf{U} = \{u_1, \dots, u_n\}$  is ένα μη εποπτευόμενο corpus από tokens,  $k$  είναι το μέγεθος του context window,  $\mathbf{P}$  μοντελοποιείται ως νευρωνικό δίκτυο με παραμέτρους  $\Theta$  (Radford et al., 2018).

# OpenAI GPT: Generative Pre-trained Transformer



Μετασχηματισμοί εισόδου για τελειοποίηση διαφορετικών εργασιών.

Μετατρέψτε όλες τις δομημένες εισόδους σε ακολουθίες διακριτικών, ακολουθούμενες linear + softmax layer.

Αρχιτεκτονική μετασχηματιστών και στόχοι εκπαίδευσης

# Ζητήματα με προηγούμενες μεθόδους

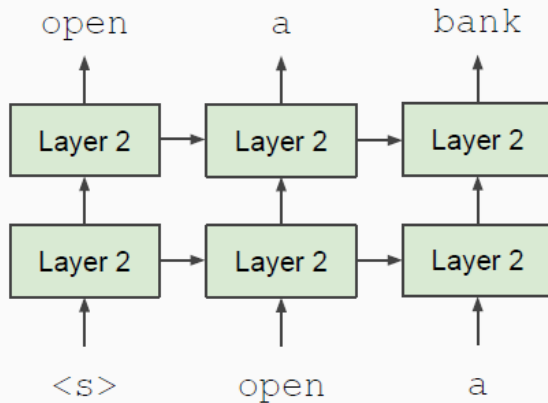
- ❑ **Πρόβλημα:** Language models Χρησιμοποιήστε μόνο αριστερό ή δεξί πλαίσιο, αλλά η κατανόηση της γλώσσας είναι αμφίδρομη.
- ❑ Γιατί τα LM είναι μονής κατεύθυνσης?
  - Λόγος 1: Η κατευθυντικότητα είναι απαραίτητη για τη δημιουργία μιας καλά διαμορφωμένης κατανομής πιθανότητας.
  - Λόγος 2: Οι λέξεις μπορούν να "δουν τον εαυτό τους" σε έναν αμφίδρομο κωδικοποιητή.



# Unidirectional vs. Bidirectional Models

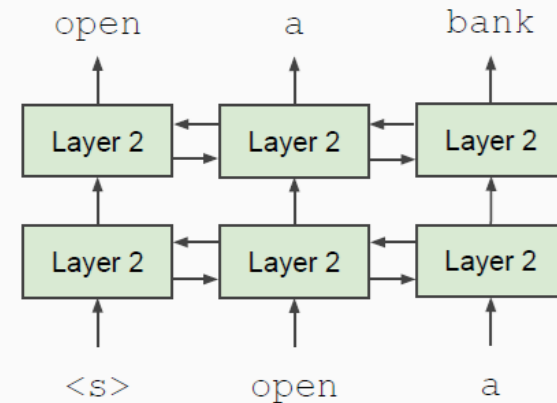
## Unidirectional context

Build representation incrementally

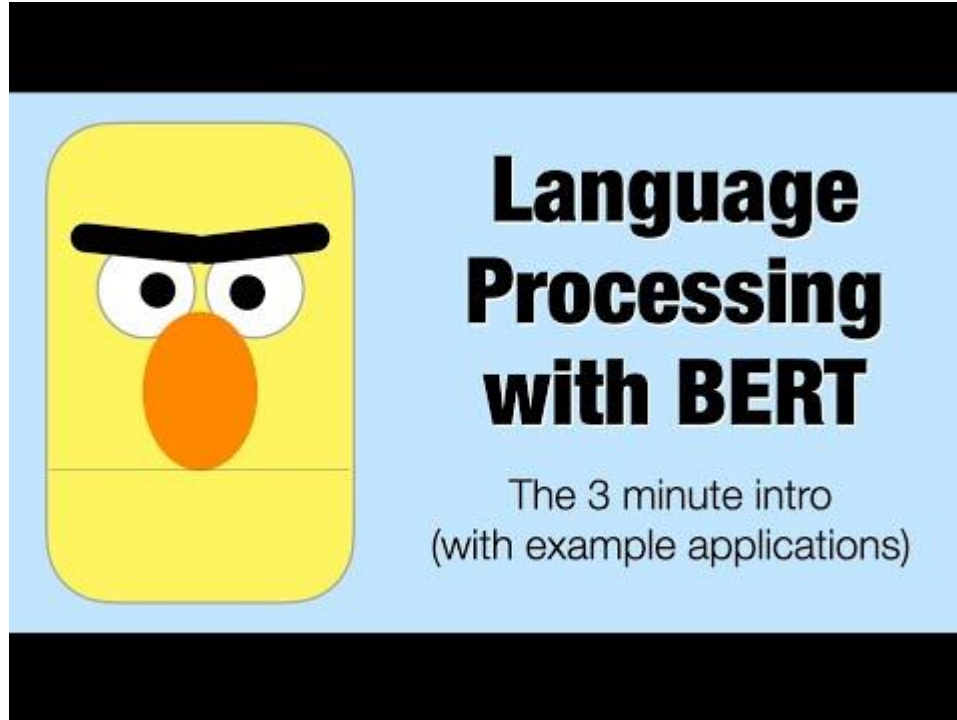


## Bidirectional context

Words can “see themselves”



# Bidirectional Encoder Representations from Transformers (BERT)



# Bidirectional Encoder Representations from Transformers (BERT)

□ To BERT εισήχθη από την Google Research το 2018 (Devlin et al., 2018)

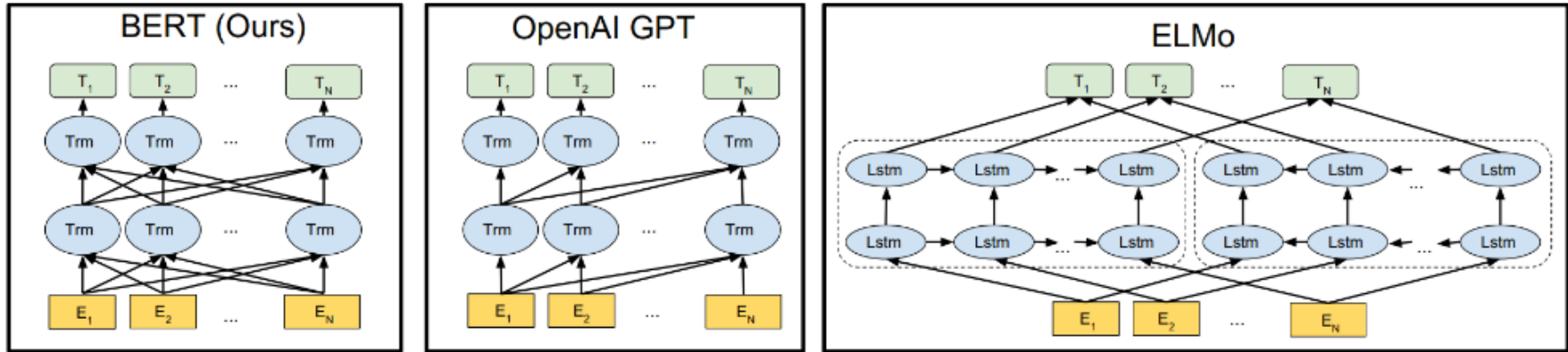
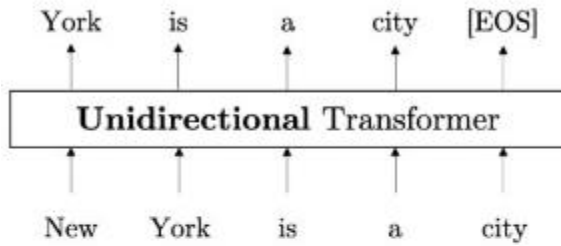


Figure 1: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks. Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

# Στόχοι για την κατάρτιση BERT

- ❑ Next-token prediction: Sequential / Unidirectional Transformer
- ❑ Masked-token prediction: Bidirectional Transformer

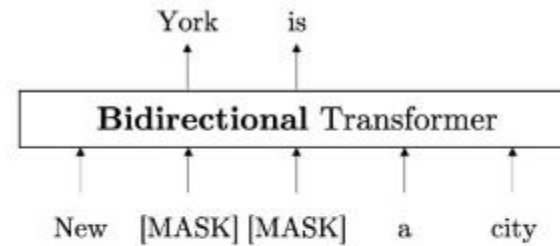
## Auto-regressive Language Modeling



$$\log p(\mathbf{x}) = \sum_{t=1}^T \log p(x_t | \mathbf{x}_{<t})$$

- Next-token prediction

## Denosing Auto-encoding (BERT)



$$\log p(\bar{\mathbf{x}} | \hat{\mathbf{x}}) = \sum_{t=1}^T \text{mask}_t \log p(x_t | \hat{\mathbf{x}})$$

- Reconstruct masked tokens

<https://github.com/zihangdai/xlnet/blob/master/misc/slides.pdf>



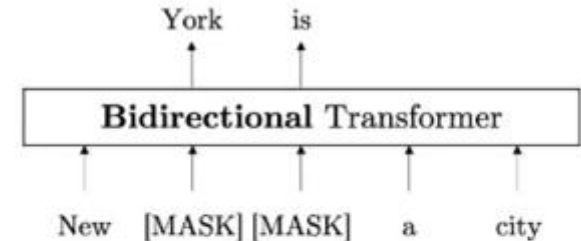
# Στόχοι για την κατάρτιση BERT

- ❑ Next-token prediction: Sequential / Unidirectional Transformer
- ❑ Masked-token prediction: Bidirectional Transformer
- Mask out  $k\%$  of the input words, and then predict the masked words.

**Input:** The man went to the **[MASK]<sub>1</sub>**.  
He bought a **[MASK]<sub>2</sub>** of milk.

**Labels:** **[MASK]<sub>1</sub>** = *store*; **[MASK]<sub>2</sub>** = *gallon*

## Denosing Auto-encoding (BERT)



$$\log p(\bar{\mathbf{x}}|\hat{\mathbf{x}}) = \sum_{t=1}^T \text{mask}_t \log p(x_t|\hat{\mathbf{x}})$$

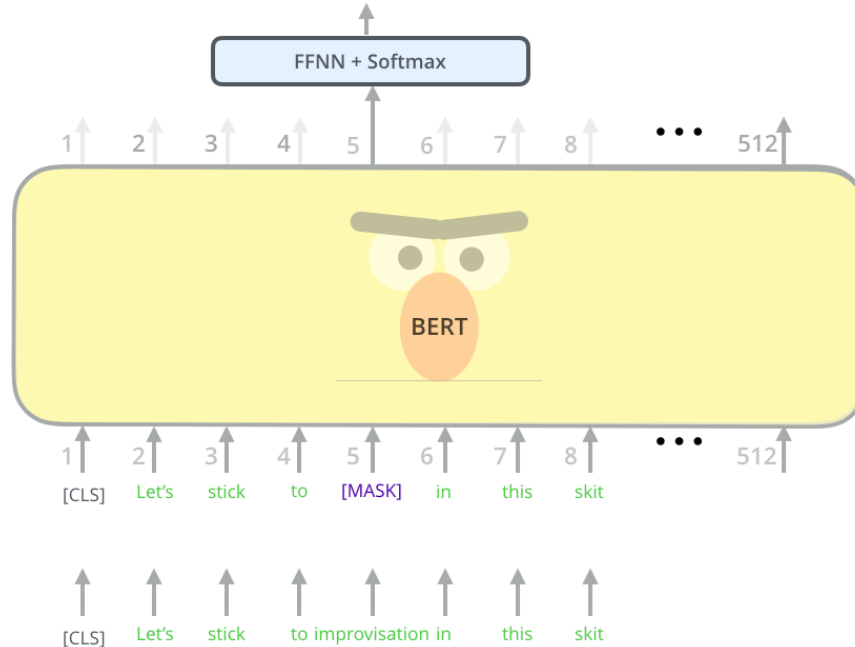
- **Reconstruct masked tokens**

# Masked Language Modeling

Use the output of the masked word's position to predict the masked word

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva



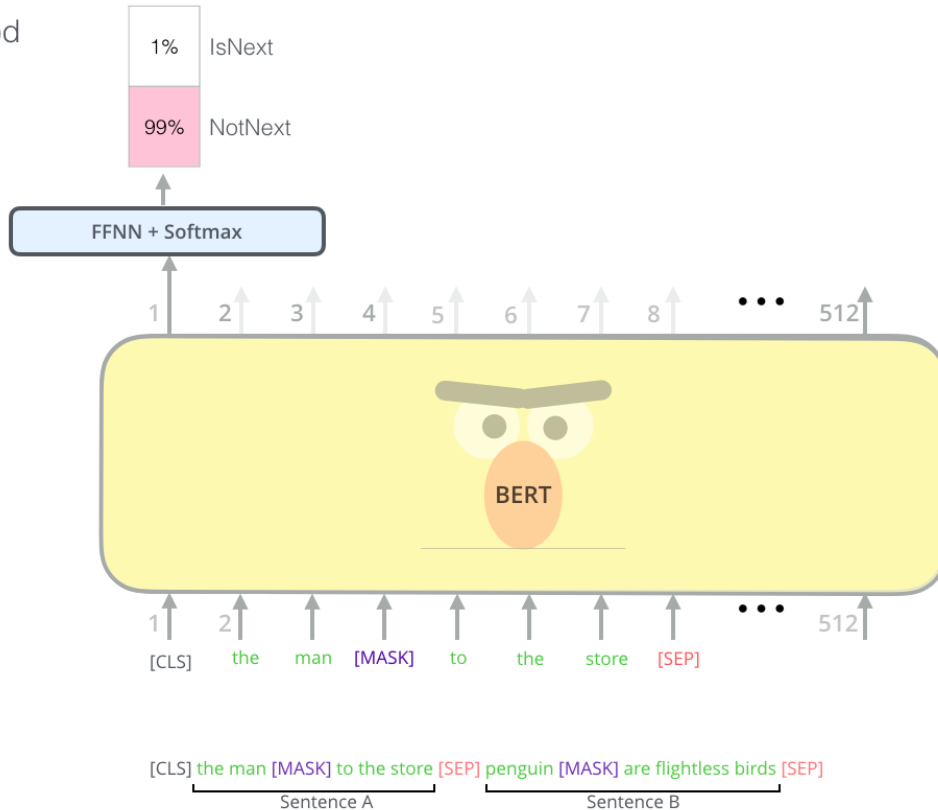
Randomly mask 15% of tokens

Input

Η έξυπνη εργασία γλωσσικής μοντελοποίησης του BERT καλύπτει το 15% των λέξεων στην είσοδο και ζητά από το μοντέλο να προβλέψει τη λέξη που λείπει.

# Next Sentence Prediction

Predict likelihood that sentence B belongs after sentence A

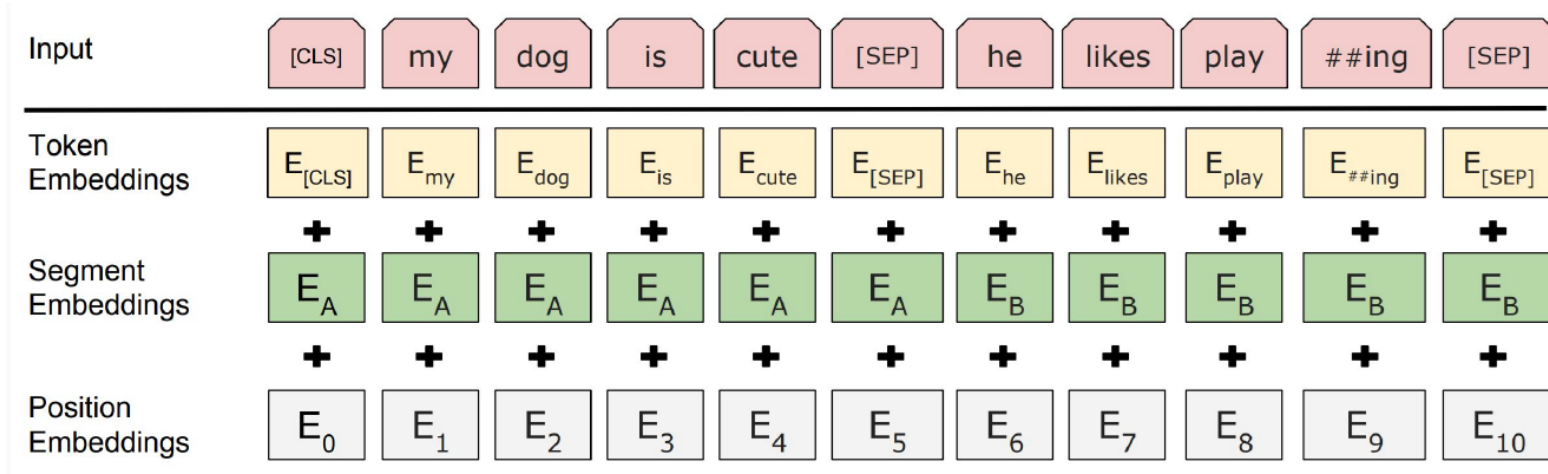


Η δεύτερη εργασία στην οποία ο BERT είναι προ-εκπαιδευμένος είναι μια εργασία ταξινόμησης δύο προτάσεων. Το tokenization υπεραπλουστεύεται σε αυτό το γραφικό, καθώς ο BERT χρησιμοποιεί στην πραγματικότητα WordPieces ως μάρκες και όχι λέξεις --- έτσι ορισμένες λέξεις αναλύονται σε μικρότερα κομμάτια.

Tokenized Input

Input

# BERT Input Representation



- ❑ Χρησιμοποιήστε 30.000 λεξιλόγιο WordPiece στην είσοδο.
- ❑ Κάθε διακριτικό είναι άθροισμα τριών ενσωματώσεων: Θέση, Τμήμα και Διακριτικό.
- ❑ Η ενιαία ακολουθία είναι πολύ πιο αποτελεσματική.

# BERT Performance Comparison

## GLUE Results

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

### MultiNLI

Premise: Hills and mountains are especially sanctified in Jainism.

Hypothesis: Jainism hates nature.

Label: Contradiction

### CoLa

Sentence: The wagon rumbled down the road.

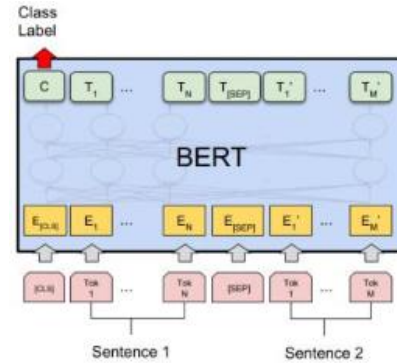
Label: Acceptable

Sentence: The car honked down the road.

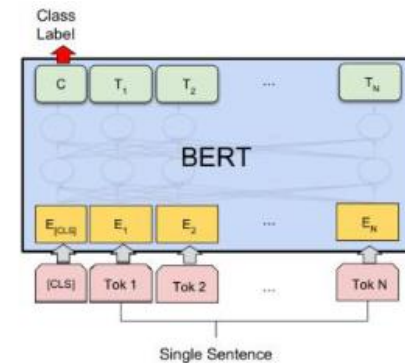
Label: Unacceptable

# Fine-Tuning BERT

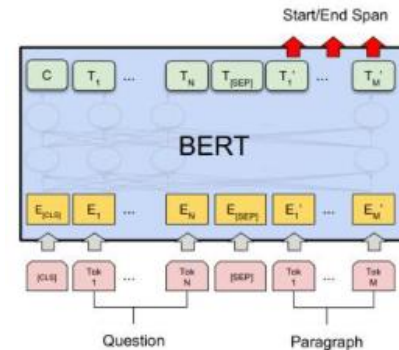
- ❑ **Context vector  $C$** : Πάρτε την τελική κρυφή κατάσταση που αντιστοιχεί στο πρώτο διακριτικό στην είσοδο [CLS].
- ❑ Μετασχηματισμός σε κατανομή πιθανότητας των ετικετών κλάσης:  $P = \text{softmax}(CW^T)$ .
- ❑ **Batch size**: 16, 32
- ❑ **Learning rate**:  $5e-5$ ,  $3e-5$ ,  $2e-5$
- ❑ **Number of epochs**: 3, 4



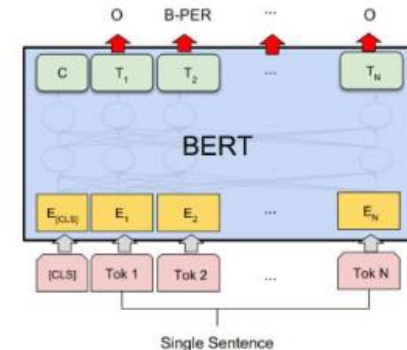
(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG



(b) Single Sentence Classification Tasks: SST-2, CoLA



(c) Question Answering Tasks: SQuAD v1.1



(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

# Περίληψη

- Εισαγωγή στις κατανεμημένες ενσωματώσεις με βάση τα συμφραζόμενα
- Ιστορία Μοντέλων Νευρωνικής Γλώσσας
  - Recurrent Neural Networks RNNs
  - ELMo
  - GPT
  - BERT
- Εσωτερική λειτουργία μετασχηματιστών και προσεγγίσεις τελειοποίησης
  - Masked Language Modeling and Next Sentence Prediction.



# Πόροι

- Jurafsky, D. and H. Martin Justin, Chapter 7. "Neural Networks and Neural Language Models" Speech and Language Processing
- Jurafsky, D. and H. Martin Justin, Chapter 11. "Transfer Learning with Contextual Embeddings and Pre-trained language models" Speech and Language Processing
- Illustrated BERT, ELMo by Jay Alammar: <https://jalammar.github.io/illustrated-bert/>
- Hugging Face's Transformers: <https://huggingface.co/docs/transformers/index>
- Transformer-based Models Hub: <https://huggingface.co/models>

