

# Επεξεργασία Φυσικής Γλώσσας & Μηχανική Μάθηση

Βαθιά Μάθηση στην ΕΦΓ

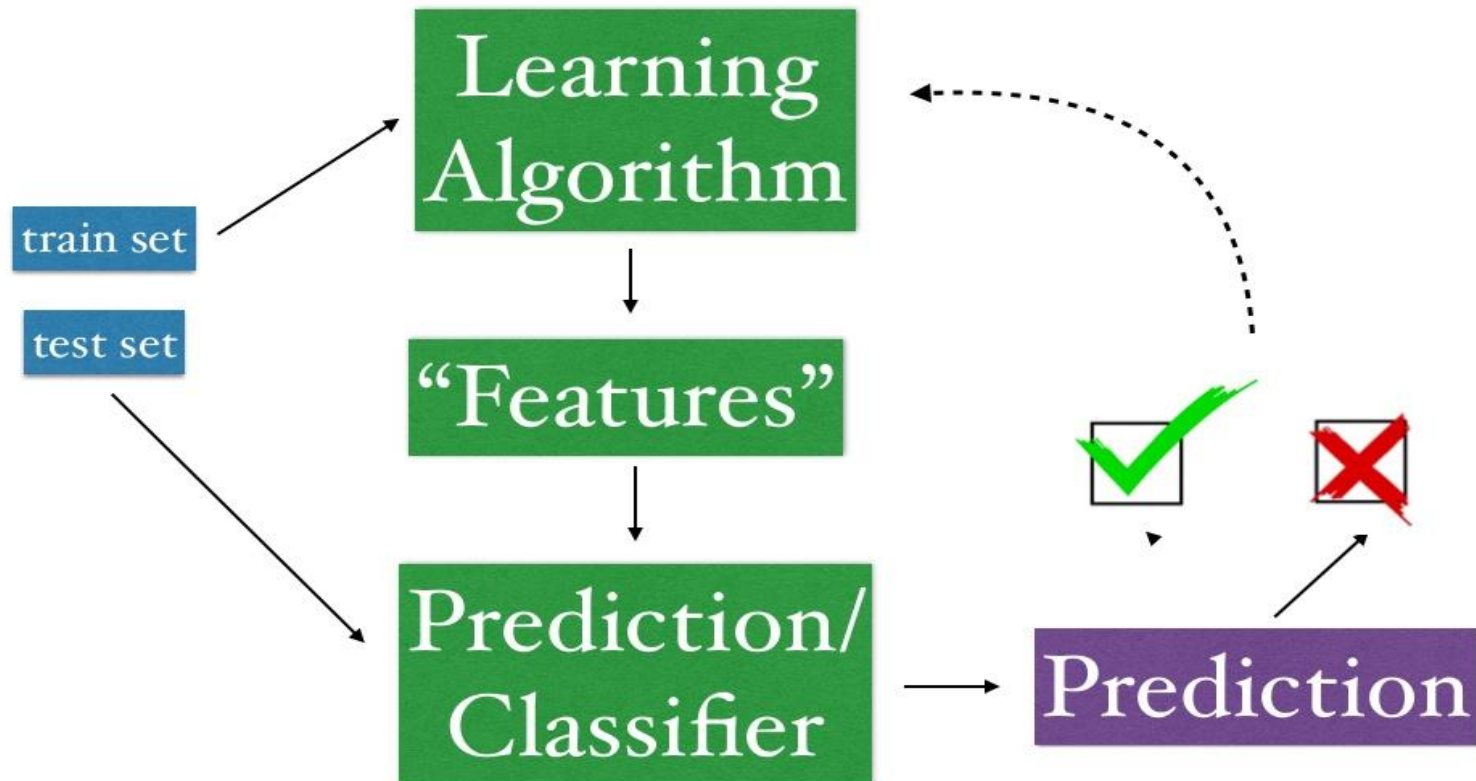
Κάτια Κερμανίδου  
kerman@ionio.gr

# Μηχανική Μάθηση: Παραδοσιακή Προσέγγιση

- Για κάθε καινούρια εργασία:
  - Σύλλεξε όσο περισσότερα επισημειωμένα δεδομένα
  - Αφιέρωσε χρόνο στην μηχανική χαρακτηριστικών (feature engineering)
    - Feature extraction
  - Τρέξε πάνω στα δεδομένα αλγορίθμους μάθησης
  - Συνέχισε την μηχανική χαρακτηριστικών
    - Feature selection
    - Dimensionality reduction
  - Επανάλαβε

# Μηχανική Μάθηση: Παραδοσιακή Προσέγγιση

## Machine Learning for NLP



# Μηχανική Μάθηση: Παραδοσιακή Προσέγγιση

- Όταν δουλεύει καλά, αυτό οφείλεται στην χειρωνακτική σχεδίαση χαρακτηριστικών για την αναπαράσταση των γλωσσολογικών φαινομένων
  - Πχ χαρακτηριστικά για την αναγνώριση ονομάτων-οντοτήτων (τοπωνύμια, ονόματα οργανισμών κλπ)
- Προβλήματα
  - Τα χειρωνακτικά σχεδιασμένα σετ χαρακτηριστικών είναι συνήθως υπερ-εξειδικευμένα, μη ολοκληρωμένα, χρονοβόρα στην σχεδίαση και την επικύρωσή τους
  - Δεν είναι αυτός ο τρόπος που μαθαίνει ο άνθρωπος.

# Πώς μαθαίνει ο άνθρωπος;

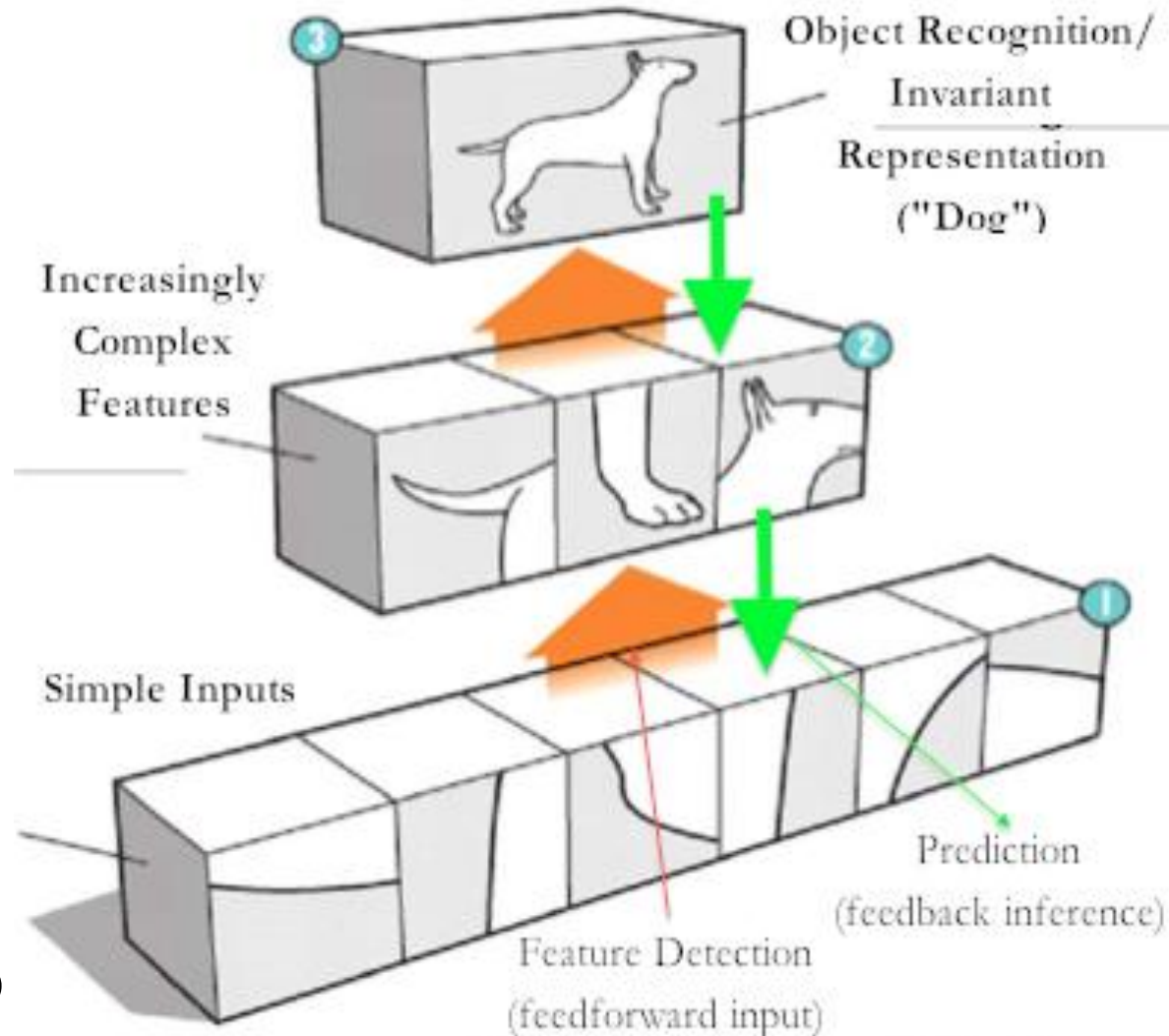
- Ο άνθρωπος μαθαίνει έννοιες και αποκτά ικανότητες και τις εφαρμόζει σε διαφορετικά προβλήματα
  - Μεταφορά (Transfer Learning)
- Ο άνθρωπος πρώτα μαθαίνει απλές έννοιες και μετά τις συνδυάζει για να μάθει πιο πολύπλοκες
- Υπάρχει η ένδειξη ότι ο πυρήνας του εγκεφάλου (cortex) έχει έναν μοναδικό αλγόριθμο μάθησης
  - Η είσοδος από τους οπτικούς νευρώνες κουναβιών δρομολογήθηκε στο τμήμα του πυρήνα που λαμβάνει την ακουστική είσοδο.
  - Είχαν την ικανότητα να μαθαίνουν να βλέπουν με αυτό το τμήμα του εγκεφάλου

# Επομένως, αν θέλουμε έναν αλγόριθμο γενικό/καθολικό, αυτός θα πρέπει να:

- Μπορεί να δουλέψει με οποιοδήποτε τύπο δεδομένων
- Μπορεί να μαθαίνει από μη επισημειωμένα δεδομένα
- Μπορεί να εξάγει αυτόματα τα χαρακτηριστικά που χρειάζεται
  - Τα χαρακτηριστικά που μαθαίνονται αυτόματα, μαθαίνονται γρήγορα και προσαρμόζονται εύκολα
- Μπορεί να μεταφέρει αυτό που έμαθε σε καινούριες εφαρμογές/περιοχές
- Μπορεί να εφαρμόζει πολυτροπική (multimodal) μάθηση
  - Να μαθαίνει ταυτόχρονα από διαφορετικές εισόδους (όραση, γλώσσα κλπ)

# Ιεραρχική Αναπαράσταση

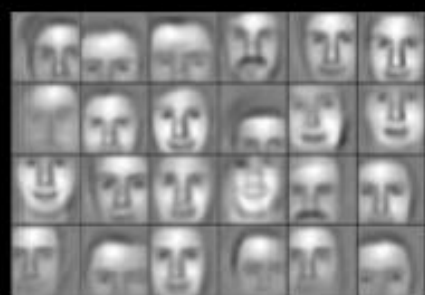
- Η αναπαράσταση του μικρόκοσμου της εφαρμογής μου γίνεται σε πολλαπλά επίπεδα
- Κάθε επίπεδο δημιουργεί καινούρια χαρακτηριστικά από συνδυασμό χαρακτηριστικών του προηγούμενου επιπέδου
- Κάθε επίπεδο είναι πιο «αφαιρετικό» (abstract) από το προηγούμενο επίπεδο



# Hierarchical Sparse coding (Sparse DBN): Trained on face images



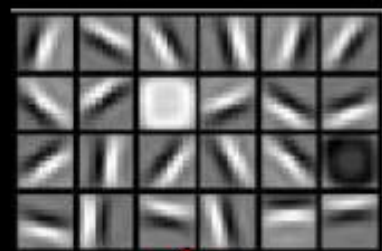
Training set: Aligned images of faces.



object models



object parts  
(combination  
of edges)



edges

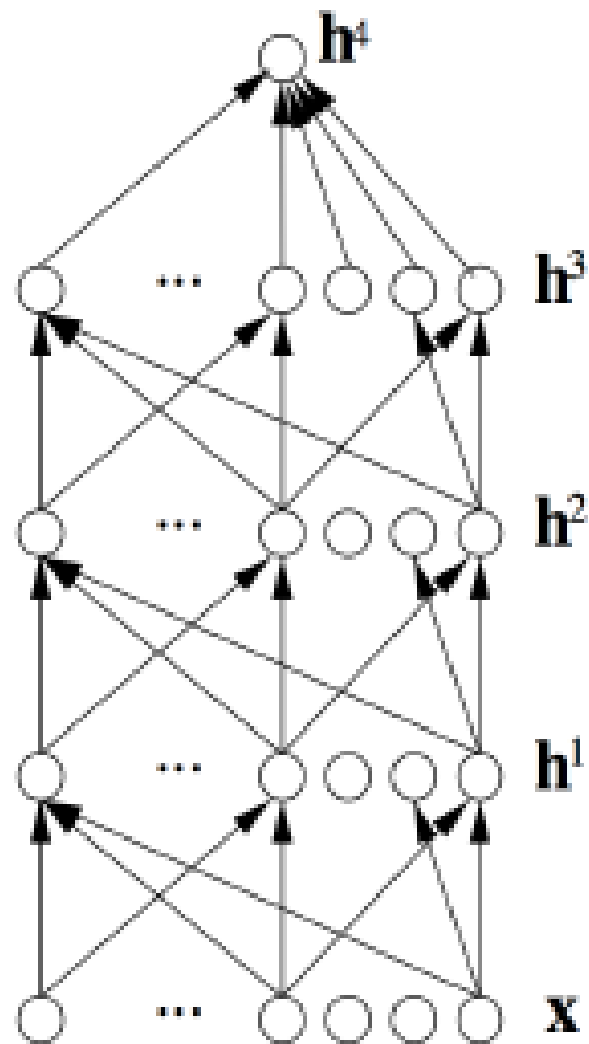


pixels



# Τι είναι Βαθιά Μάθηση (Deep Learning):

- Οι αλγόριθμοι βαθιάς μάθησης στοχεύουν να μάθουν πολλαπλά επίπεδα αναπαράστασης (χαρακτηριστικών) και μια έξοδο
- Από ένα σύνολο ωμών (raw) δεδομένων εισόδου  $x$  (πχ λέξεις)
- Η βασική οικογένεια αλγορίθμων είναι τα νευρωνικά δίκτυα



# Γιατί Βαθιά Μάθηση;

- Οι αλγόριθμοι βαθιάς μάθησης
  - είναι και επιβλεπόμενοι και μη-επιβλεπόμενοι
  - Παρέχουν ένα καθολικό πλαίσιο για την πολυεπίπεδη αναπαράσταση γνώσης
  - Το 2006 ξεκίνησαν να αποδίδουν καλύτερα από την παραδοσιακή μάθηση
  - Σήμερα
    - Έχουν την δυνατότητα πρόσβασης σε περισσότερα δεδομένα
    - Έχουν πρόσβαση σε μεγαλύτερη υπολογιστική ισχύ
    - Στηρίζονται σε νέους αλγορίθμους και αρχιτεκτονικές

# Deep Learning: Why for NLP ?

One Model rules them all ?

DL approaches have been successfully applied to:

Automatic summarization

Coreference resolution

Discourse analysis

Machine translation

Morphological segmentation

Named entity recognition (NER)

Natural language generation

Word sense disambiguation

Relationship extraction

Speech processing

Part-of-speech tagging

sentence boundary disambiguation

Sentiment analysis

Optical character recognition (OCR)

Question answering

Parsing

Word segmentation

Natural language understanding

Information retrieval (IR)

Speech recognition

Topic segmentation and recognition

Speech segmentation

Information extraction (IE)

# Βαθιά Μάθηση στην ΕΦΓ

- Μεγάλες βελτιώσεις τα τελευταία χρόνια
  - Στα διάφορα επίπεδα γλωσσολογικής πληροφορίας
    - φωνητικό
    - Μορφολογικό
    - Συντακτικό
    - Σημασιολογικό
  - Στις διάφορες εφαρμογές
    - Αυτόματη μετάφραση
    - Ανάλυση συναισθήματος
    - Συστήματα ερωταποκρίσεων

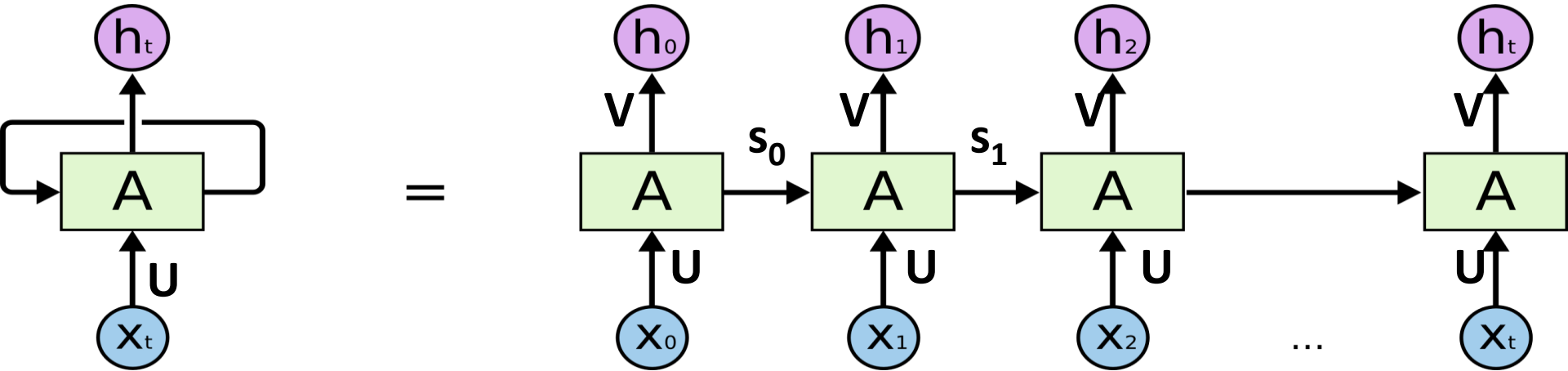
# Βαθιές Αρχιτεκτονικές

- Είδη νευρωνικών δικτύων
  - Recursive neural networks (RNN)
    - Recurrent neural networks
    - Long short term memory neural networks (LSTM)
  - Convolutional neural networks (CNN)
  - Sequence-to-sequence models

# Φαινόμενα ακολουθίας στην γλώσσα (Sequential data)

- Για την αποσαφήνιση του νοήματος σε προτάσεις, λέξεις, χαρακτήρες χρειάζονται τα συμφραζόμενά τους.
- Μηχανική Μετάφραση
  - Μια λέξη έχει διαφορετικό νόημα ανάλογα με τα συμφραζόμενά της
- Ανάλυση Συναισθήματος
  - Η εμφάνιση επιρρημάτων και λέξεων άρνησης (όπως "very", "not", και "a bit too") στα συμφραζόμενα της λέξης που κρύβει το συναίσθημα επηρεάζουν την ένταση, την πόλωση ή την αντιστροφή του συναισθήματος.
- Διαλογικά συστήματα
  - Το επόμενο βήμα σε έναν διάλογο καθορίζεται από τα προηγούμενα βήματα του διαλόγου και τον στόχο που έχει ο διάλογος.
- Tokenization
  - Οι προηγούμενοι και οι επόμενοι χαρακτήρες χρησιμοποιούνται για να αναγνωριστεί η έναρξη μιας καινούριας λέξης.

# Recurrent NNS



Ένα RNN είναι μια αλυσίδα αντιγράφων του ίδιου δικτύου. Οι ίδιες συνάψεις, τα ίδια βάρη, εφαρμόζονται σε καινούρια είσοδο (πχ καινούρια λέξη) σε κάθε χρονικό βήμα.

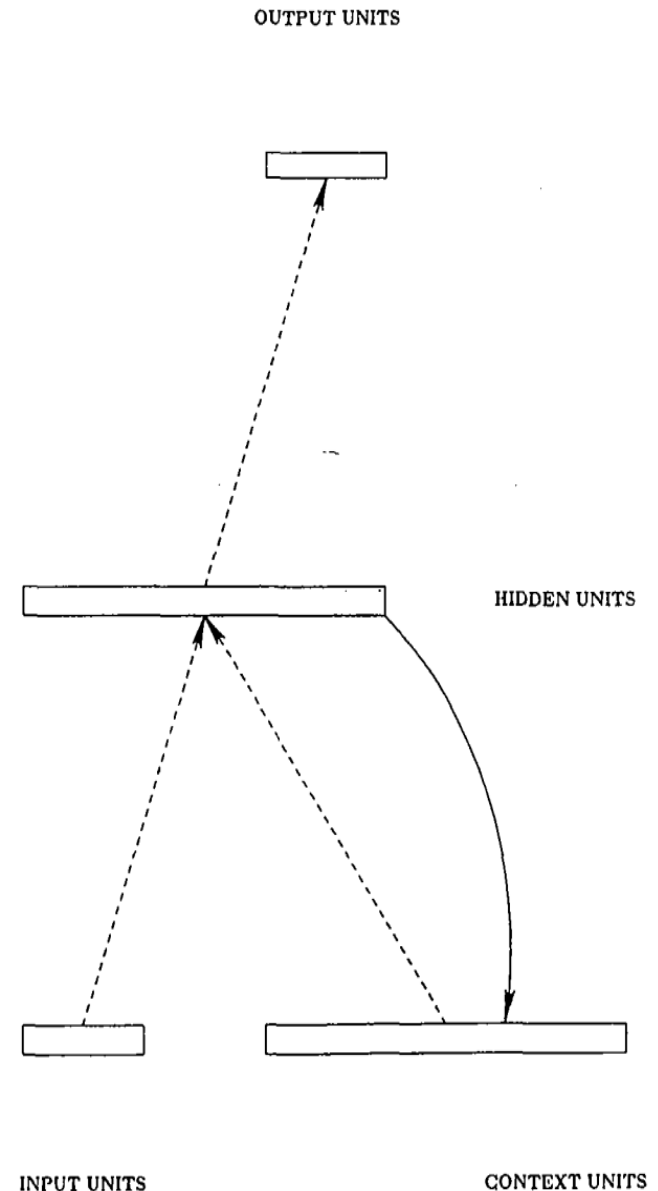
Τα RNNs συνδυάζουν την τρέχουσα είσοδο με την κατάσταση του προηγούμενου βήματος σε μια συνάρτηση η οποία παράγει την καινούρια τρέχουσα κατάσταση.

$$s_t = f(Ux_t + Ws_{t-1})$$

$$h_t = \text{softmax}(Vs_t)$$

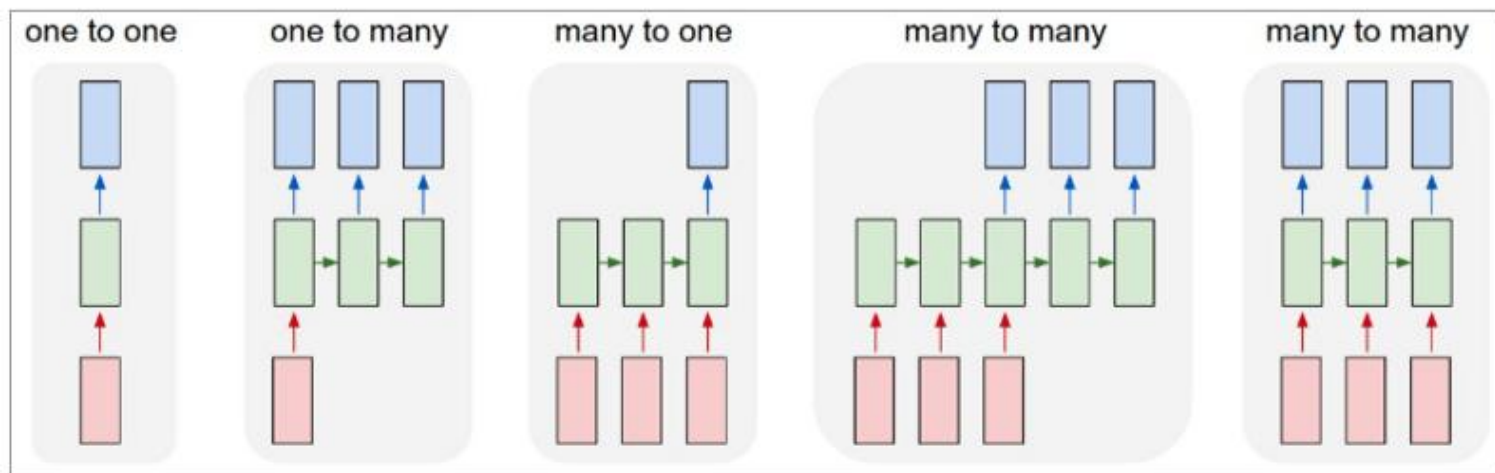
# Recurrent NNS

- Το πλεονέκτημα των RNN είναι η ικανότητά τους να αντιμετωπίζουν ακολουθιακά δεδομένα, χάρη στη «μνήμη» τους. Ενώ τα νευρωνικά δίκτυα δεν έχουν αίσθηση του χρόνου, και η πρόβλεψή τους εξαρτάται από την τωρινή τους είσοδο μόνο, τα RNNs λαμβάνουν υπόψη τους και την τωρινή είσοδο και την «είσοδο συμπραζομένων» ("context unit"), η οποία «χτίζεται» βάσει των όσων έχουν δει προηγούμενα.
- Έτσι, η πρόβλεψη που πραγματοποιείται την στιγμή  $T$  επηρεάζεται από αυτήν που πραγματοποιήθηκε την στιγμή  $T-1$ .





# ΑΡΧΙΤΕΚ ΤΟΝΙΚΕΣ RNN



- Vanilla mode of processing without RNN, from fixed-sized input to fixed-sized output (e.g. image classification).
- Sequence output (e.g. image captioning takes an image and outputs a sentence of words).
- Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment)
- Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French).
- Synced sequence input and output (e.g. video classification where we wish to label each frame of the video).

# Trigram RNN via POS tagging

Here,  $h^{(t)}$  not only depends on the previous hidden state  $h^{(t-1)}$ , but also directly depends on  $h^{(t-2)}$ . We hope that this extra dependency can help to catch longer windows in the sentence.

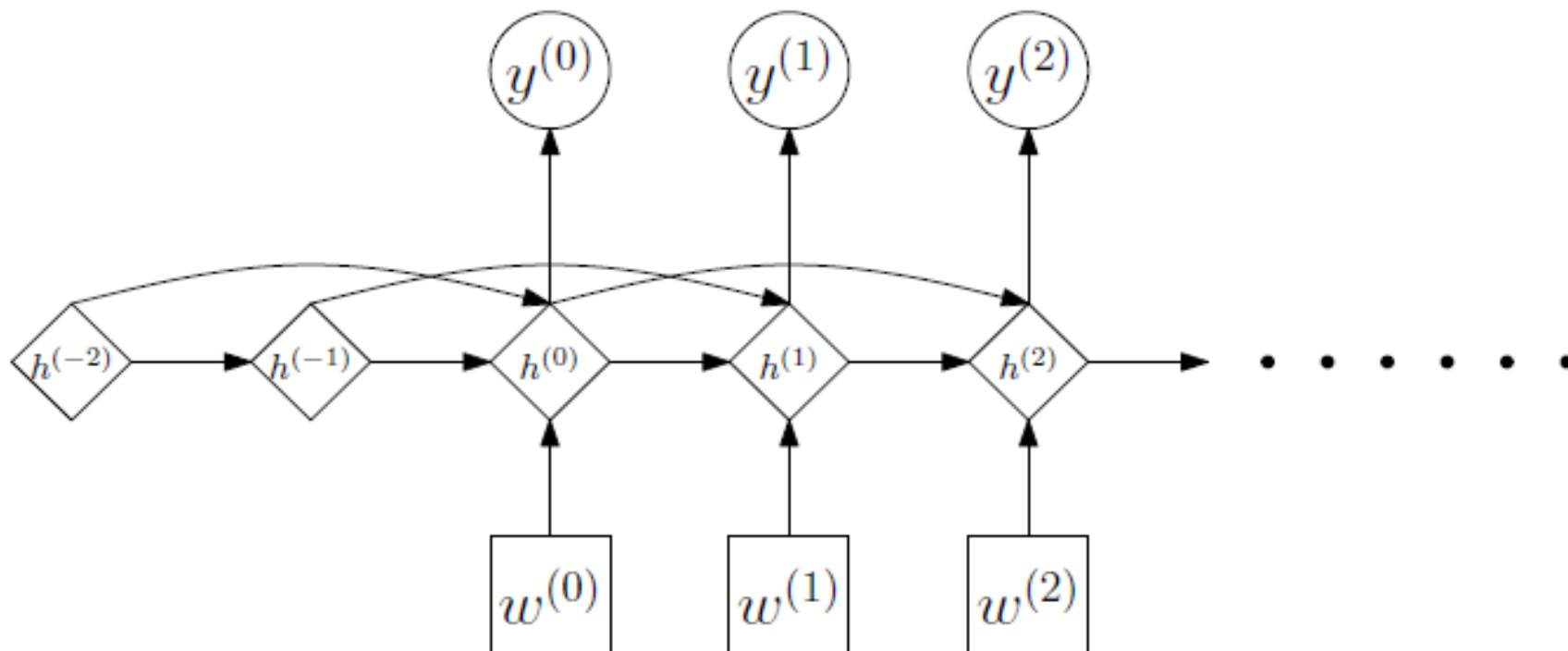


Figure 2: Trigram RNN

- <https://cs224d.stanford.edu/reports/QinLonglu.pdf>

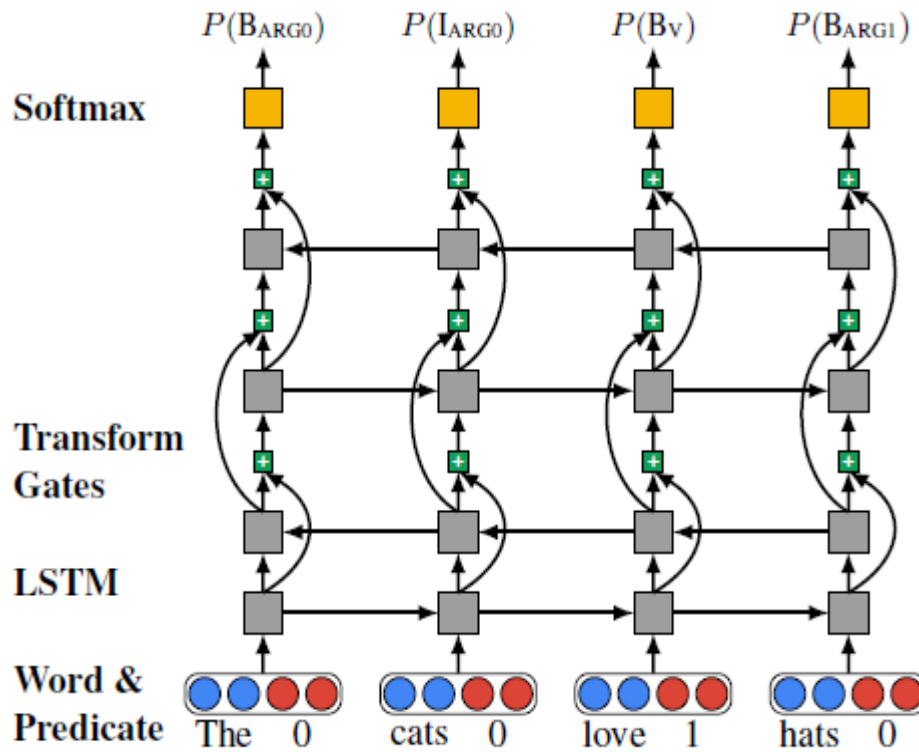
# RNNS για Semantic Role Labeling

## Deep Semantic Role Labeling: What Works and What's Next, He et al., 2017

Ρόλος	Περιγραφή	Παράδειγμα
<b>AO</b>	Agent, ο υποκινητής της ενέργειας του ρήματος	<b>The workers</b> dumped large sacks
<b>A1</b>	Theme, αυτό που δέχεται την ενέργεια του ρήματος	The workers dumped <b>large sacks</b>
<b>A2</b>	Έμμεσο, δεύτερο αντικείμενο	He added 5 <b>to the initial calculation</b>
<b>A3</b>	Beneficiary	He gave John 300€ <b>for the car.</b>
<b>DIR</b>	κατεύθυνση	The workers dumped large sacks <b>into a huge bin</b>
<b>ADV</b>	Επιρρηματικός προσδιορισμός	<b>Maybe</b> he is sick.
<b>LOC</b>	Τοπικός προσδιορισμός	He studied <b>in Germany.</b>
<b>MNR</b>	Τροπικός προσδιορισμός	The workers <b>mechanically</b> dumped large sacks
<b>PNC</b>	σκοπός	He saves money <b>to pay for the car.</b>

# RNNS για Semantic Role Labeling

## Deep Semantic Role Labeling: What Works and What's Next, He et al., 2017



Έξοδος: το tag του ρόλου της λέξης

B-Beginning

$B_{ARGO}$  - αρχή της φράσης του AO

I-Inside

$I_{ARGO}$  - εσωτερικό της φράσης του AO

O-Outside (εκτός φράσης)

Είσοδος: το διάνυσμα της λέξης, ακολουθούμενο από 1 αν είναι το ρήμα, 0 αλλιώς.

Figure 1: Highway LSTM with four layers. The curved connections represent highway connections, and the plus symbols represent transform gates that control inter-layer information flow.

# Confusion matrix

pred. \ gold	A0	A1	A2	A3	ADV	DIR	LOC	MNR	PNC	TMP
A0	-	55	11	13	4	0	0	0	0	0
A1	78	-	46	0	0	22	11	10	25	14
A2	11	23	-	48	15	56	33	41	25	0
A3	3	2	2	-	4	0	0	0	25	14
ADV	0	0	0	4	-	0	15	29	25	36
DIR	0	0	5	4	0	-	11	2	0	0
LOC	5	9	12	0	4	0	-	10	0	14
MNR	3	0	12	26	33	0	0	-	0	21
PNC	0	3	5	4	0	11	4	2	-	0
TMP	0	8	5	0	41	11	26	6	0	-

Table 5: Confusion matrix for labeling errors, showing the percentage of predicted labels for each gold label. We only count predicted arguments that match gold span boundaries.

# Μέτρο f για διάφορες αποστάσεις ρήματος-ρόλου

(απόσταση = αριθμός λέξεων που παρεμβάλλονται)

- Pradhan, Punyakanok: ρηχά μοντέλα
- L2: δυο κρυφά επίπεδα
- L4: τέσσερα κρυφά επίπεδα
- L6: έξι κρυφά επίπεδα
- L8: οκτώ κρυφά επίπεδα

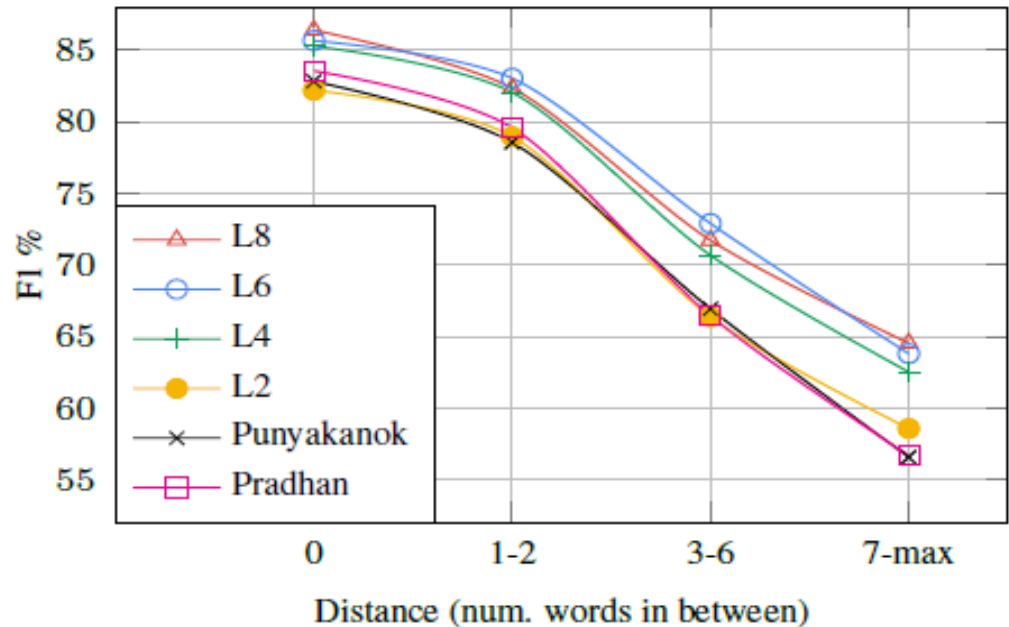
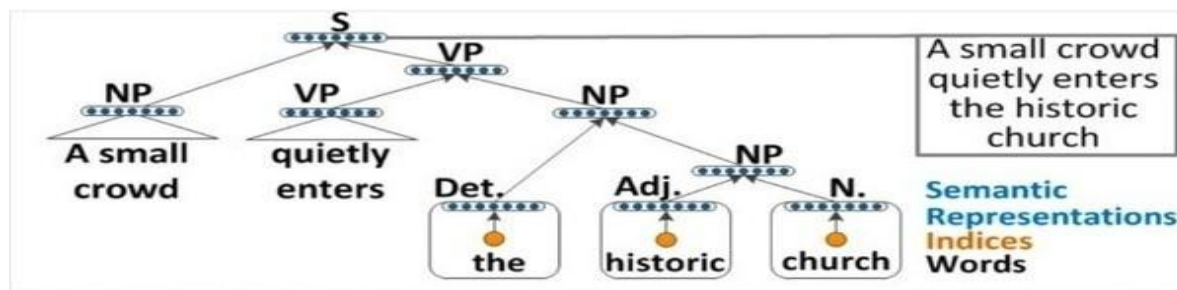


Figure 5: F1 by surface distance between predicates and arguments. Performance degrades least rapidly on long-range arguments for the deeper neural models.

# Recursive NNS

- Στα Recursive NNS δεν έχω την έννοια των χρονικών βημάτων.
- Είναι ιεραρχικά δίκτυα στα οποία η είσοδος πρέπει να υποστεί ιεραρχική επεξεργασία σε μορφή δενδρικής δομής.
- Στην παρακάτω εικόνα φαίνεται πώς ένα recursive NN μαθαίνει το συντακτικό δέντρο μιας πρότασης παίρνοντας αναδρομικά την έξοδο της πράξης που πραγματοποιήθηκε σε ένα μικρότερο κομμάτι του κειμένου.

## Recursive Neural Tensor Network



# Long short term memory NNS

- Στα Recurrent NNs το να λάβω υπόψη πολλά χρονικά βήματα πριν μπορεί να
  - Προκαλέσει *exploding gradients* (αύξηση βαρών πολύ απότομη λόγω του επαναλαμβανόμενου πολλαπλασιασμού των βαρών που είναι μεγαλύτερα της μονάδας)
  - Προκαλέσει *vanishing gradients* (μείωση βαρών πολύ απότομη λόγω του επαναλαμβανόμενου πολλαπλασιασμού των βαρών που είναι μικρότερα της μονάδας)
- Για αυτό έχουν προταθεί τα Long short term memory NNS
- Είναι RNNs που περιλαμβάνουν ένα κύτταρο (LSTM unit)
- Επιτρέπουν στα RNNs να μαθαίνουν για πολλά χρονικά βήματα (πάνω από 1000).
- Το κύτταρο ρυθμίζει την διέλευση της πληροφορίας μέσα στο δίκτυο. Αποφασίζει τι θα αποθηκεύσει, τι θα διαβαστεί, τι θα διαγραφεί μέσω πυλών που ανοιγοκλείνουν.



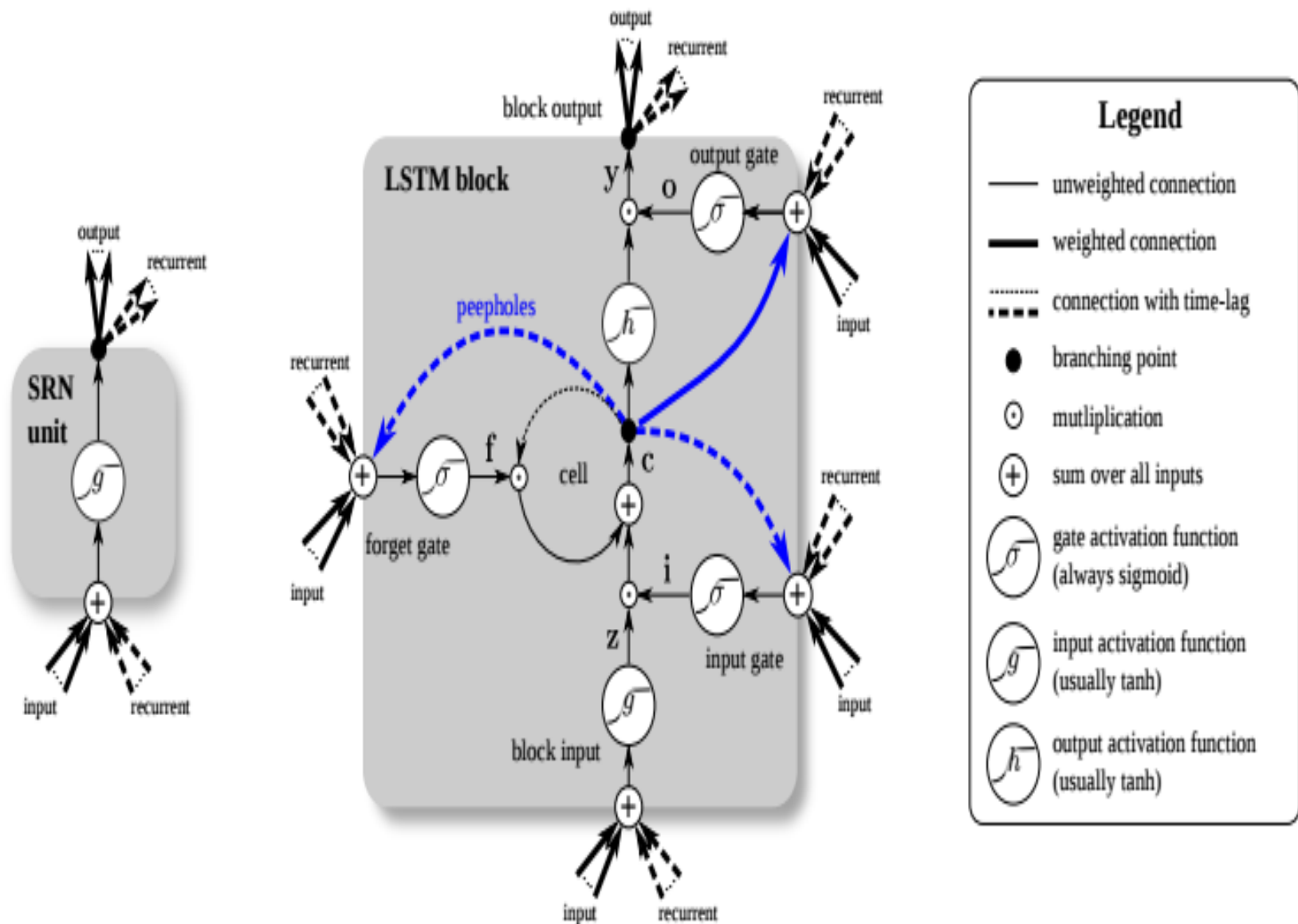


Figure 1. Detailed schematic of the Simple Recurrent Network (SRN) unit (left) and a Long Short-Term Memory block (right) as used in the hidden layers of a recurrent neural network.

# Convolution - Συνέλιξη

To sliding window

1	0	1
0	1	0
1	0	1

ονομάζεται πυρήνας (kernel),

φίλτρο (filter), ή ανιχνευτής χαρακτηριστικών (*feature detector*). Εδώ χρησιμοποιείται ένα φίλτρο 3×3, οι τιμές του πολλαπλασιάζονται κελί-κελί με τον αρχικό πίνακα, και αθροίζονται. Για την πλήρη συνέλιξη πραγματοποιείται αυτό για κάθε στοιχείο, «τσουλώνοντας» το φίλτρο πάνω από όλον τον αρχικό πίνακα.

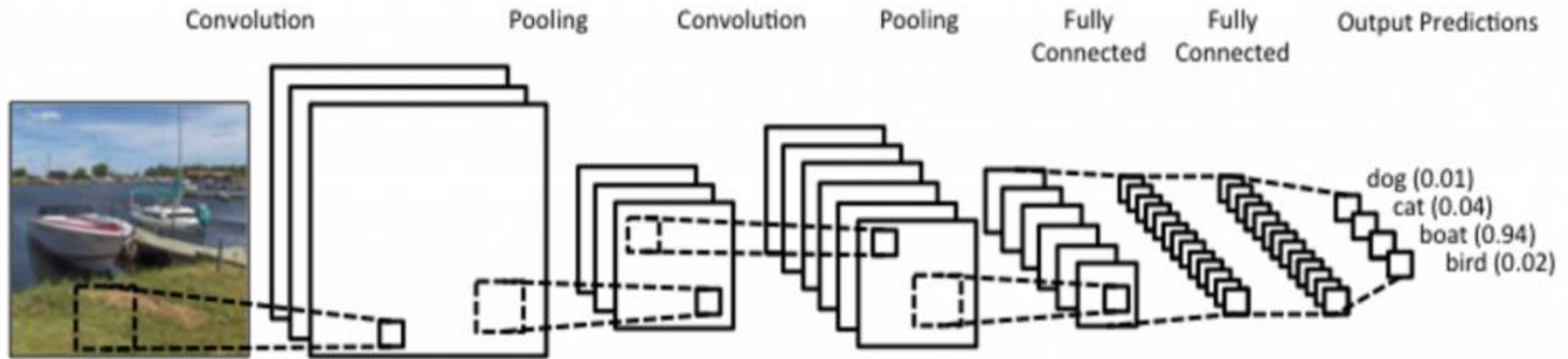
1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved  
Feature

# Convolutional NNs



Σε ένα convolutional layer (CL) δεν έχω πλήρεις συνδέσεις, αλλά πραγματοποιώ συνελίξεις πάνω στο επίπεδο εισόδου και η έξοδος του CL είναι το αποτέλεσμα της συνέλιξης.

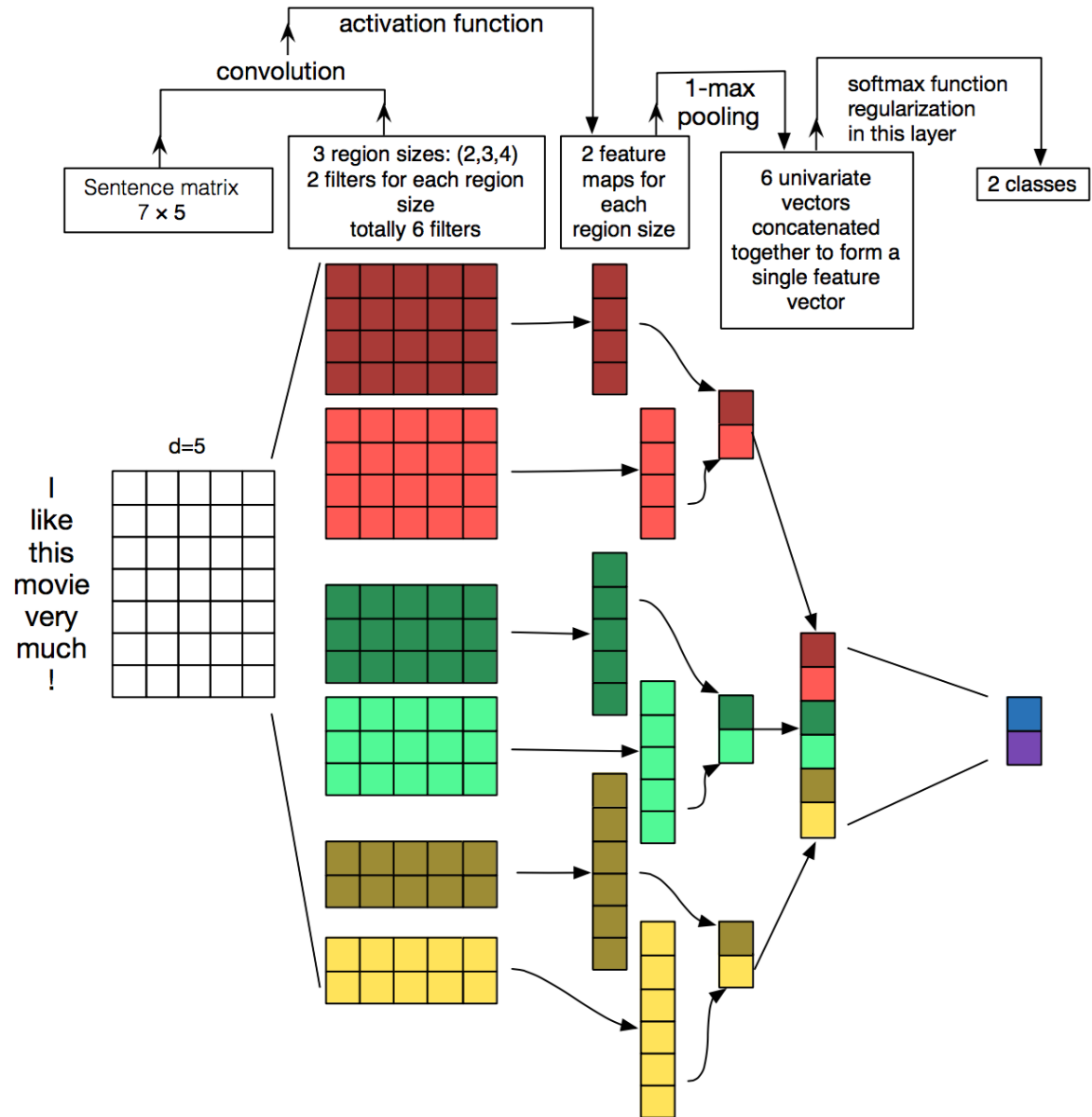
Δημιουργούνται μόνο τοπικές συνδέσεις, όπου κάθε περιοχή της εισόδου συνδέεται με έναν νευρώνα της εξόδου.

Κατά την εκπαίδευση μαθαίνονται οι τιμές του φίλτρου.

Για παράδειγμα, στην ταξινόμηση εικόνων, σε ένα πρώτο CL μαθαίνεται η αναγνώριση ακμών, σε δεύτερο CL χρησιμοποιούνται οι ακμές για να μάθει το δίκτυο απλά σχήματα, και σε τρίτο CL χρησιμοποιούνται τα σχήματα για να μάθει το δίκτυο πιο υψηλού επιπέδου στοιχεία της εικόνας, όπως πχ σχήματα ανθρώπινου προσώπου.

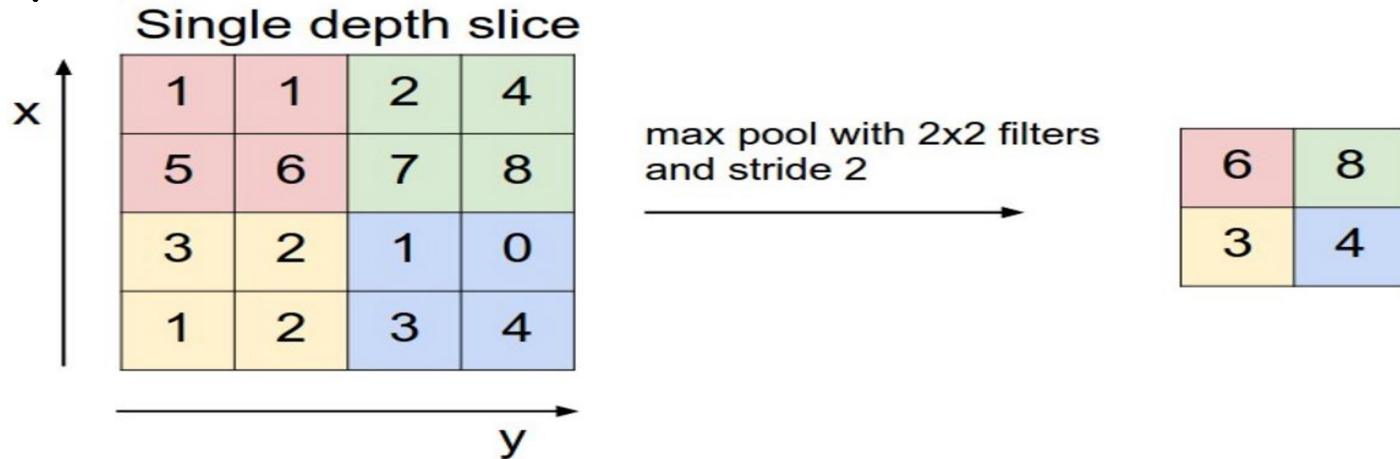
# Convolutional NNs στην ΕΦΓ

Εδώ ο αρχικός μου πίνακας δεν είναι τιμές pixel, αλλά γραμμές από word embeddings. Κάθε γραμμή είναι το διάνυσμα μιας λέξης σε μια πρόταση/κείμενο.



# Convolutional NNs στην ΕΦΓ - Pooling

Max pooling for a 2x2 window (in NLP we typically apply pooling over the complete output, yielding just a single number for each filter).



Με το pooling κάνω υποδειγματοληψία της εξόδου του CL. Με αυτό πετυχαίνω:

- Σταθερού μεγέθους έξοδο (αν πχ θέλω ταξινόμηση σε κλάση δύο τιμών μπορώ να ρυθμίσω της έξοδο ώστε να είναι διάνυσμα δύο θέσεων), ανεξάρτητα από το μέγεθος της εισόδου (το μήκος των προτάσεών μου)
- Να μειώνω την διαστατικότητα μου, διατηρώντας την σημαντική πληροφορία

# Convolutional Neural Network for Paraphrase Identification, Yin & Schuetze, 2015

Με είσοδο δυο προτάσεων , βγάζω δυαδική έξοδο αν έχουν το ίδιο νόημα ή όχι.

- (1) "Mary gave birth to a son in 2000."
- (2) "He is 19 years old and his mother is Mary."

Κάθε επίπεδο του συνελικτικού ΝΔ αναπαριστά και πιο πολύπλοκες γλωσσολογικές οντότητες:

- μονόγραμμα
- Μικρά n-γραμμα
- Μεγάλα n-γραμμα
- πρόταση

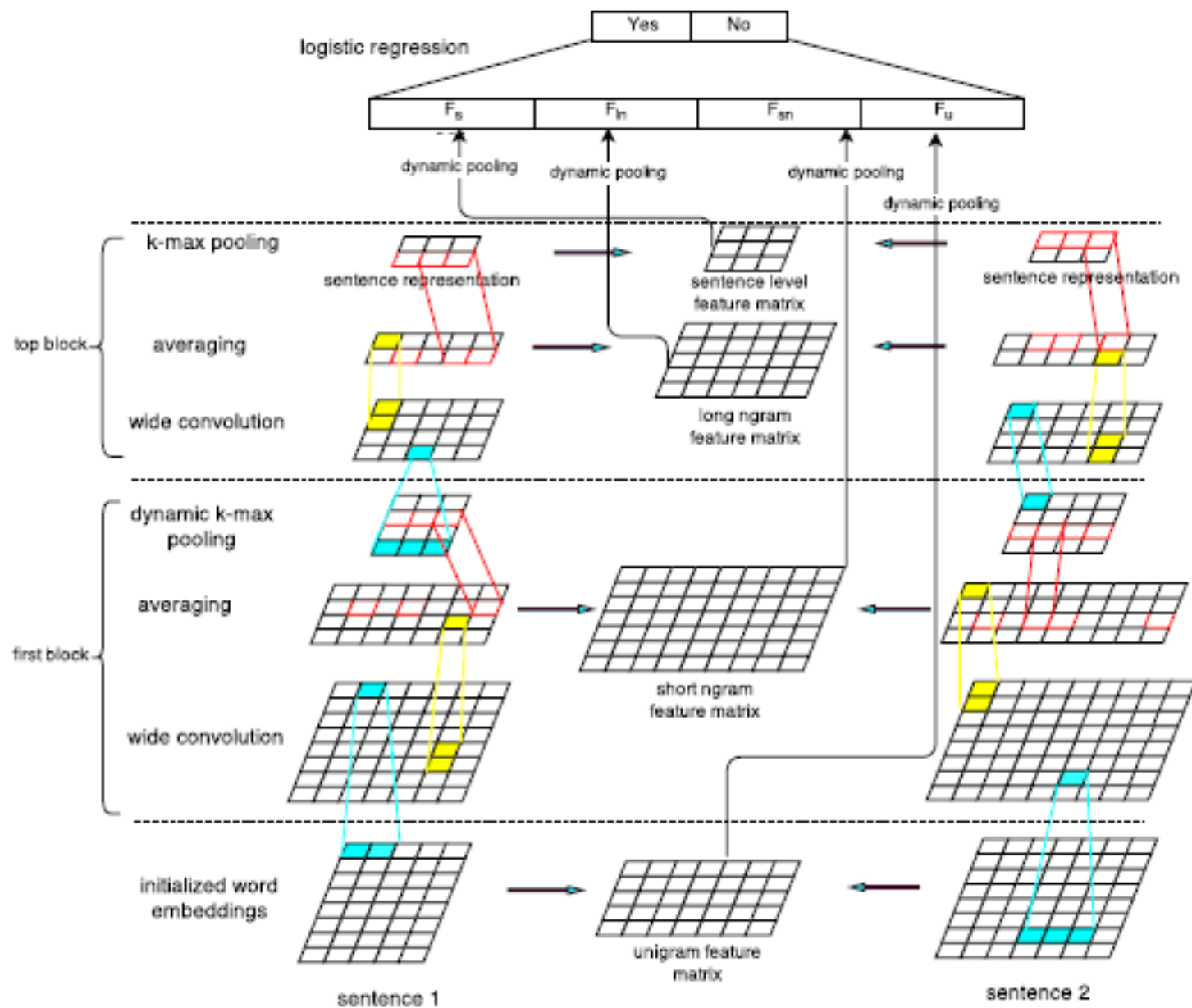


Figure 1: The paraphrase identification architecture Bi-CNN-MI

	features used	acc	$F_1$
1	no features	66.5	79.9
2	+ u: unigram	68.4	79.7
3	+ sn: short ngram	75.3	82.8
4	+ ln: long ngram	76.2	83.1
5	+ s: sentence	73.4	82.3
6	– u: unigram	77.8	84.3
7	– sn: short ngram	76.3	83.5
8	– ln: long ngram	75.6	83.2
9	– s: sentence	77.6	84.2
10	all features	78.1	84.4

Table 2: Analysis of impact of the four feature classes. Line 1: majority baseline. Line 10: Bi-CNN-MI result from Table 1. Lines 2–5: Bi-CNN-MI when only one feature class is used. Line 6–9: ablation experiment: on each line one feature class is removed.



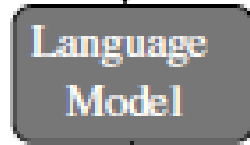
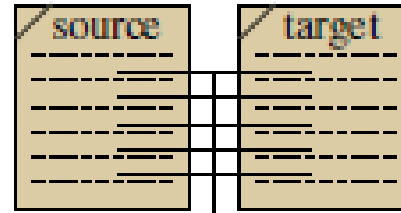
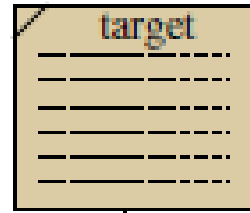
# Στατιστική Μηχανική Μετάφραση

## SMT Architecture

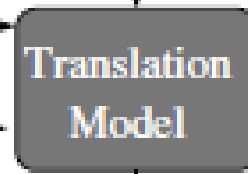
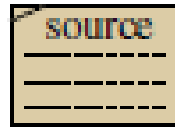
- Πολλά δομικά στοιχεία
- Αρχιτεκτονική pipeline
- Μετάφραση ανεξάρτητη συμφραζομένων
- Ανεξάρτητα μοντέλα (γλώσσας, μετάφρασης)
- Κάθε παράμετρος βελτιστοποιείται τοπικά

n-gram extraction

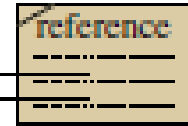
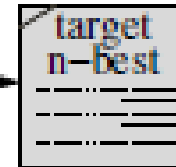
word alignment



*training*



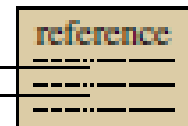
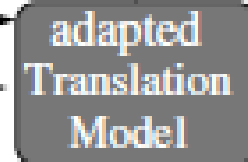
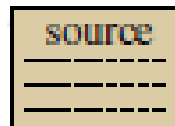
update



converged

scoring

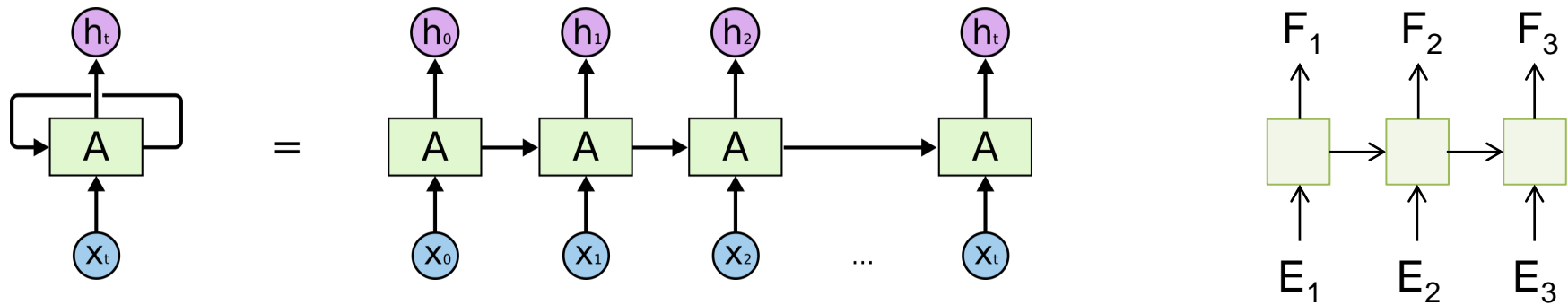
*tuning*



scoring

*testing*

# RNN - Μηχανική Μετάφραση

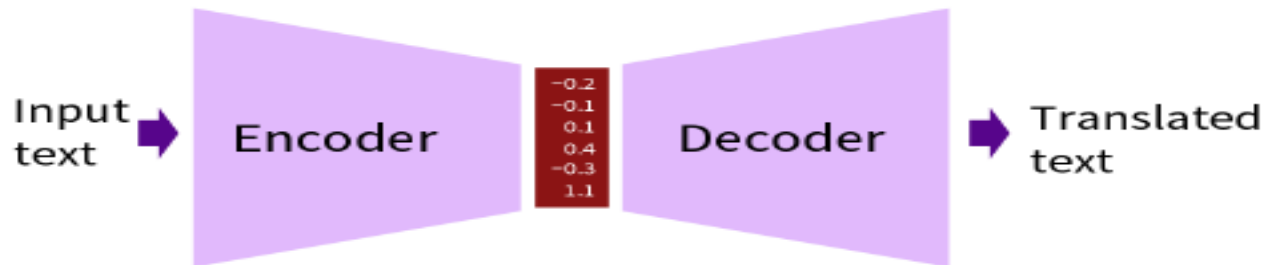


Μια RNN αρχιτεκτονική σαν τις παραπάνω δεν είναι κατάλληλη για μετάφραση μιας πρότασης από μια γλώσσα σε μια άλλη γιατί

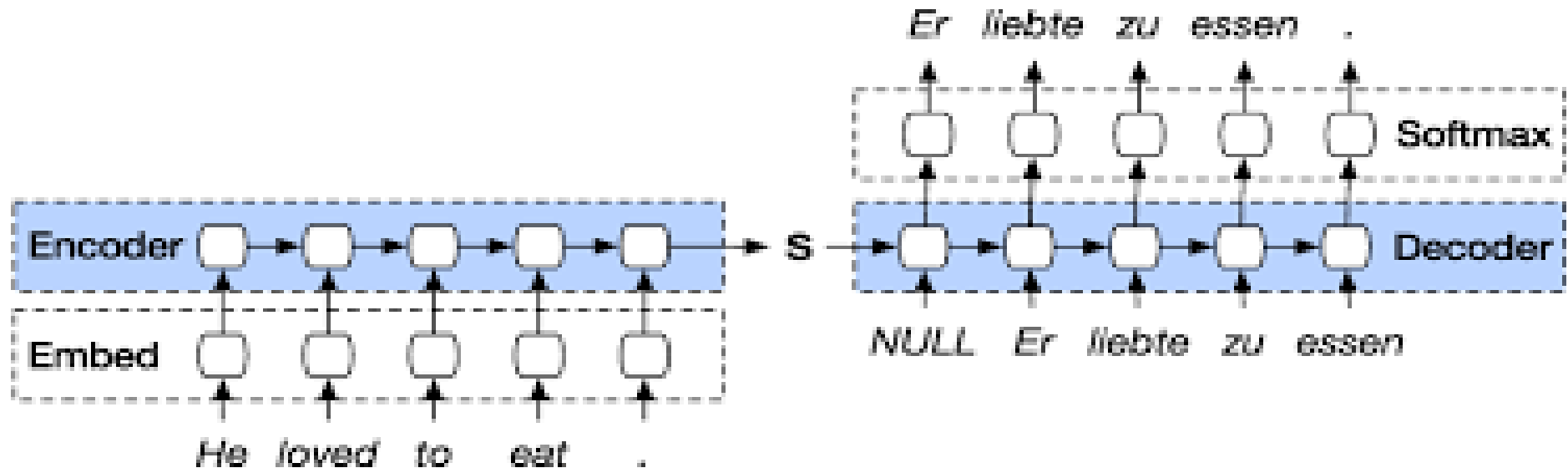
- Οι δυο προτάσεις μπορεί να είναι διαφορετικού μήκους
- Οι λέξεις πολύ πιθανό να μην είναι ευθυγραμμισμένες

Οπότε πρέπει να διαβαστεί ολόκληρη η πρόταση πηγή προτού μεταφραστεί.

# Encoder-Decoder Αρχιτεκτονικές



- Όλα σε ένα μεγάλο μοντέλο
- Καθολική βελτιστοποίηση των παραμέτρων
- Δεν υπάρχει ρητή κατάτμηση σε διακριτά μοντέλα



# Μοντέλα Seq2Seq - Εκπαίδευση

Για Sequence to sequence modeling έχουν προταθεί encoder-decoder αρχιτεκτονικές:

**Εκπαίδευση (παράλληλα σώματα κειμένων):**

Ο encoder (RNN) παίρνει σαν είσοδο τα embeddings της πρότασης πηγής

Ο decoder (RNN) παίρνει σαν είσοδο τα embeddings της πρότασης στόχου

Τα embeddings αναπαρίστανται μέσω κρυφών καταστάσεων

- Η έξοδος κάθε κρυφής κατάστασης εξαρτάται
  - από την τωρινή είσοδο
  - από την έξοδο της προηγούμενης κρυφής κατάστασης

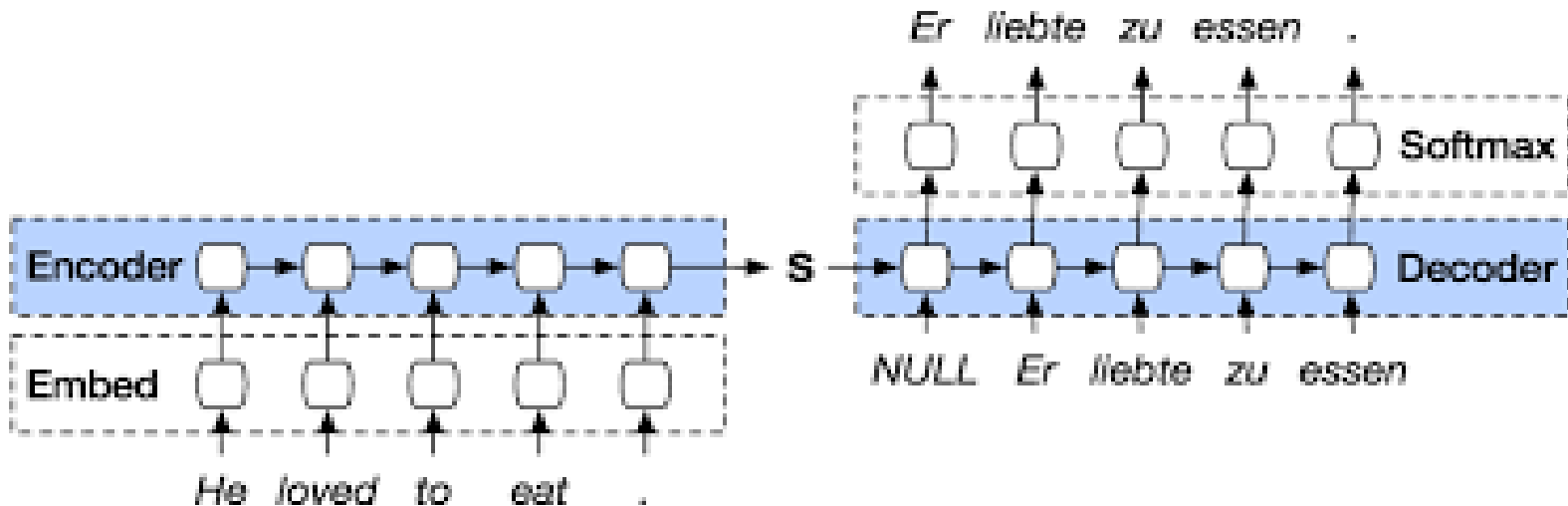
Η αρχική κατάσταση του encoder είναι τυχαία.

Η αρχική κατάσταση του decoder είναι η τελευταία κρυφή κατάσταση του encoder ( $s$ )

Κατά το backpropagation «μαθαίνονται» τα βάρη στον encoder, προκειμένου να μάθει καλύτερες διανυσματικές αναπαραστάσεις για τις προτάσεις και ταυτόχρονα «μαθαίνονται» τα βάρη του decoder για να μάθει να παράγει γραμματικά σωστές προτάσεις που είναι σχετικές με το διάνυσμα συμφραζομένων  $s$ .

# Μοντέλα Seq2Seq - Μετάφραση

- Ο encoder θα πάρει την πρόταση πηγή, και θα βγάλει το  $s$ .
- Ο decoder θα ενεργοποιηθεί με το που θα δει το σύμβολο NULL (EOS).
- Στον πρώτο του κρυφό κόμβο θα πάρει σαν είσοδο το διάνυσμα  $s$ .
- Διατρέχονται όλα τα επίπεδα, εφαρμόζεται η  $\text{softmax}()$  στη έξοδο του τελευταίου επιπέδου και έτσι προβλέπεται η πρώτη λέξη εξόδου.
- Η λέξη αυτή επανατροφοδοτείται σαν είσοδος στο NN, μαζί με το διάνυσμα  $s$ , διατρέχονται πάλι όλα τα επίπεδα, εφαρμόζεται η  $\text{softmax}()$  στη έξοδο του τελευταίου επιπέδου και έτσι προβλέπεται η δεύτερη λέξη εξόδου.
- Με τον ίδιο τρόπο προβλέπονται και οι υπόλοιπες λέξεις της εξόδου.



# Βιβλιογραφία/ Δικτυογραφία

- <http://cs224d.stanford.edu/>
- <http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>
- <https://recast.ai/blog/ml-spotlight-rnn/>
- <https://skymind.ai/wiki/lstm>
- <https://adeshpande3.github.io/A-Beginner%27s-Guide-To-Understanding-Convolutional-Neural-Networks/>
- <http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/>
- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>, 2015
- [https://web.stanford.edu/class/cs224n/archive/WW\\_WW\\_1617/lecture\\_notes/cs224n-2017-notes6.pdf](https://web.stanford.edu/class/cs224n/archive/WW_WW_1617/lecture_notes/cs224n-2017-notes6.pdf)