

# A Unified Architecture for Natural Language Processing

(Deep Neural Networks with Multitask Learning)

**Ronan Collobert**

ronan@collobert.com

&

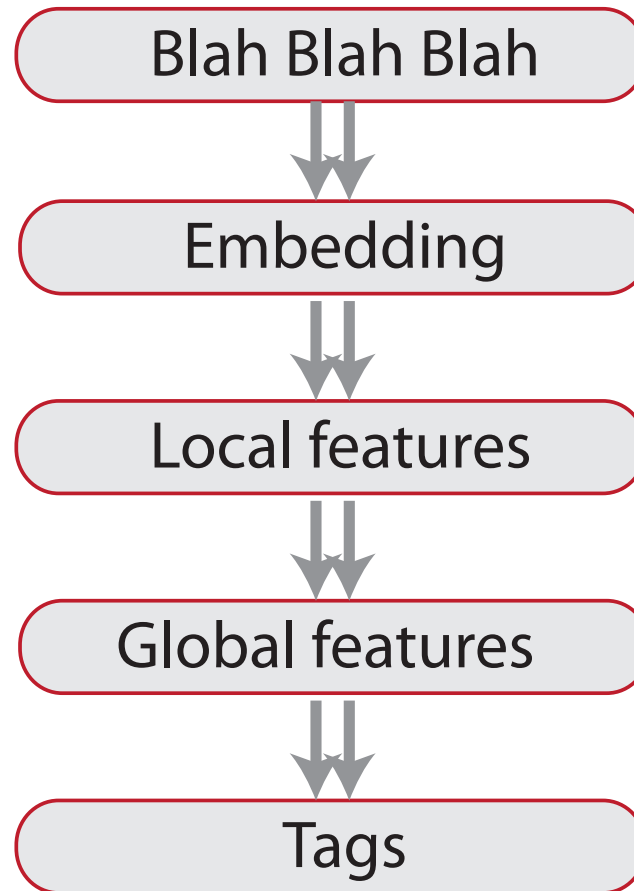
**Jason Weston**

jaseweston@gmail.com

NEC Laboratories America

# The Big Picture

(1/2)



**Deep** architecture



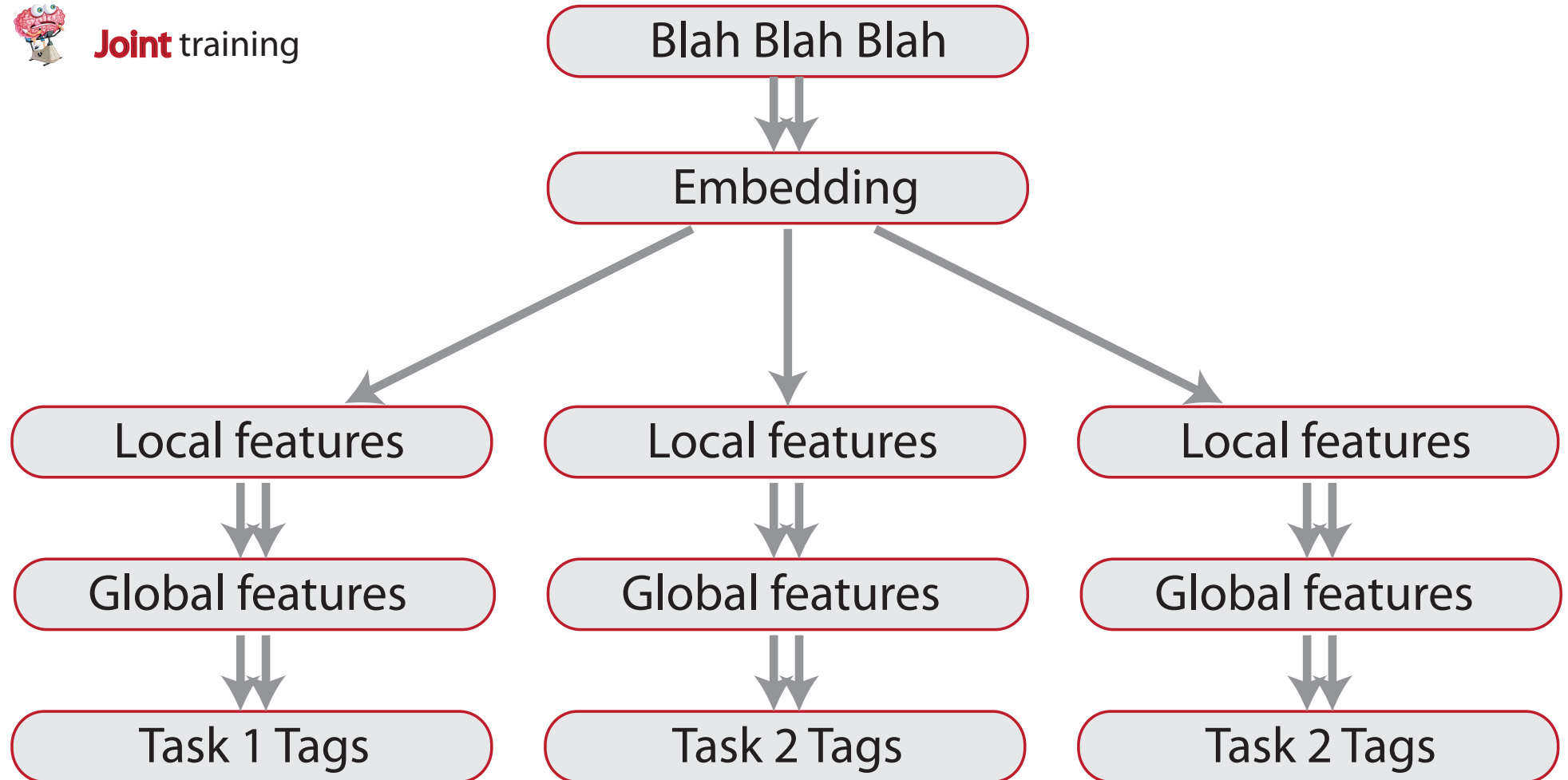
**Unification** of NLP tasks

# The Big Picture

(2/2)



**Joint** training



# NLP Tasks

---



Part-Of-Speech Tagging (POS): syntactic roles (noun, adverb...)



Chunking: syntactic constituents (noun phrase, verb phrase...)



Name Entity Recognition (NER): person/company/location...



Semantic Role Labeling (SRL): semantic role

[John]*ARG0* [ate]*REL* [the apple]*ARG1* [in the garden]*ARGM-LOC*

---

Labeled data: Wall Street Journal ( $\sim 1M$  of words)

# The Shallow System Way

(1/2)



Choose some good **hand designed features**

<p><b>Predicate and POS tag</b> of predicate</p> <p><b>Phrase type:</b> adverbial phrase, prepositional phrase, . . .</p> <p><b>Head word</b> and POS tag of the head word</p> <p><b>Path:</b> traversal from predicate to constituent</p> <p><b>Word-sense</b> disambiguation of the verb</p> <p><b>Length</b> of the target constituent (number of words)</p> <p><b>Partial Path:</b> lowest common ancestor in path</p> <p><b>First and last words</b> and POS in constituents</p> <p><b>Constituent tree distance</b></p> <p><b>Dynamic class context:</b> previous node labels</p> <p><b>Constituent relative features:</b> head word</p> <p><b>Constituent relative features:</b> siblings</p>	<p><b>Voice:</b> active or passive (hand-built rules)</p> <p><b>Governing category:</b> Parent node's phrase type(s)</p> <p><b>Position:</b> left or right of verb</p> <p>Predicted <b>named entity</b> class</p> <p><b>Verb clustering</b></p> <p><b>NEG</b> feature: whether the verb chunk has a "not"</p> <p><b>Head word replacement</b> in prepositional phrases</p> <p><b>Ordinal position</b> from predicate + constituent type</p> <p><b>Temporal cue words</b> (hand-built rules)</p> <p><b>Constituent relative features:</b> phrase type</p> <p><b>Constituent relative features:</b> head word POS</p> <p><b>Number of pirates existing in the world. . .</b></p>
--	--



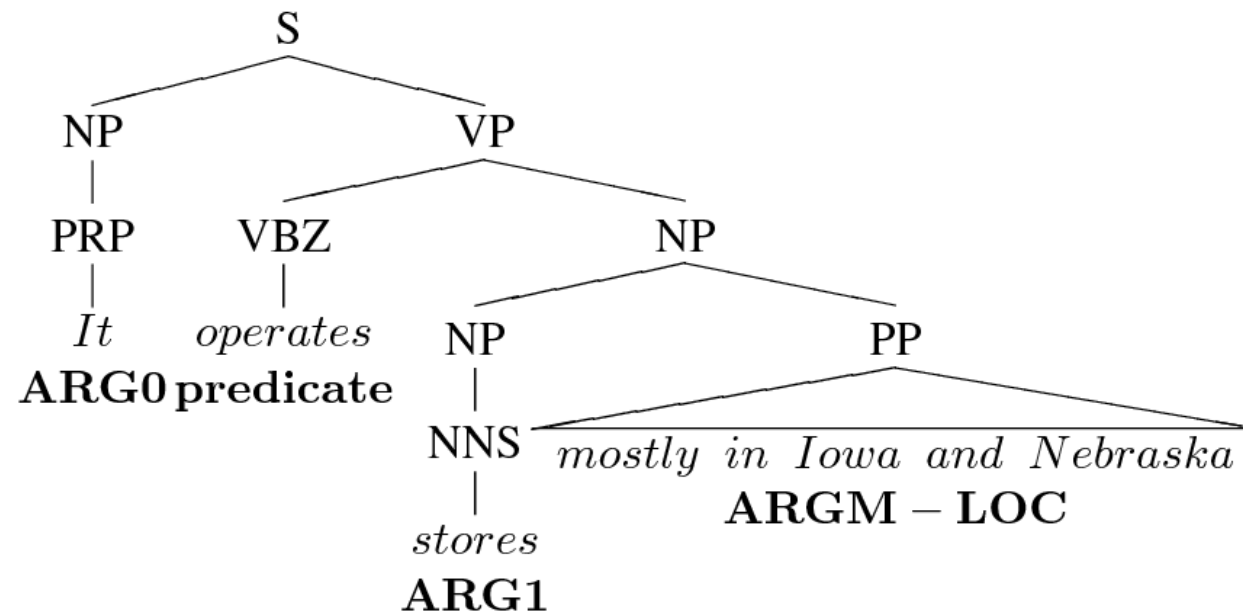
Feed them to a **shallow classifier** like SVM

# The Shallow System Way

(2/2)



Cascade features: e.g. extract POS, construct a parse tree



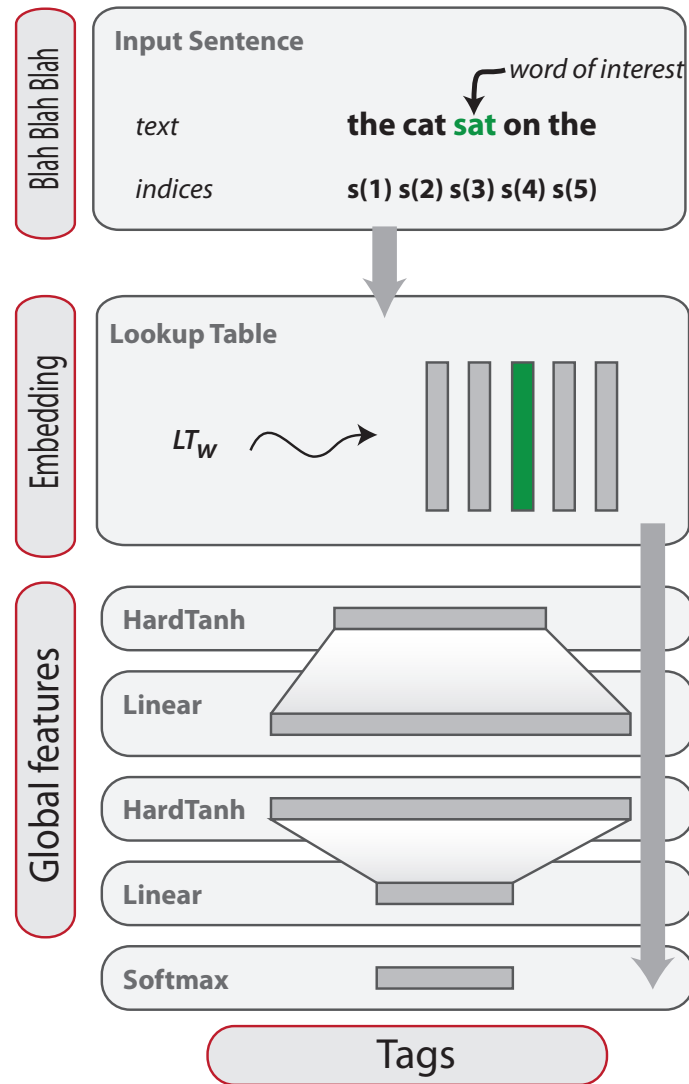
Extract hand-made features from the parse tree



Feed these features to a shallow classifier like SVM

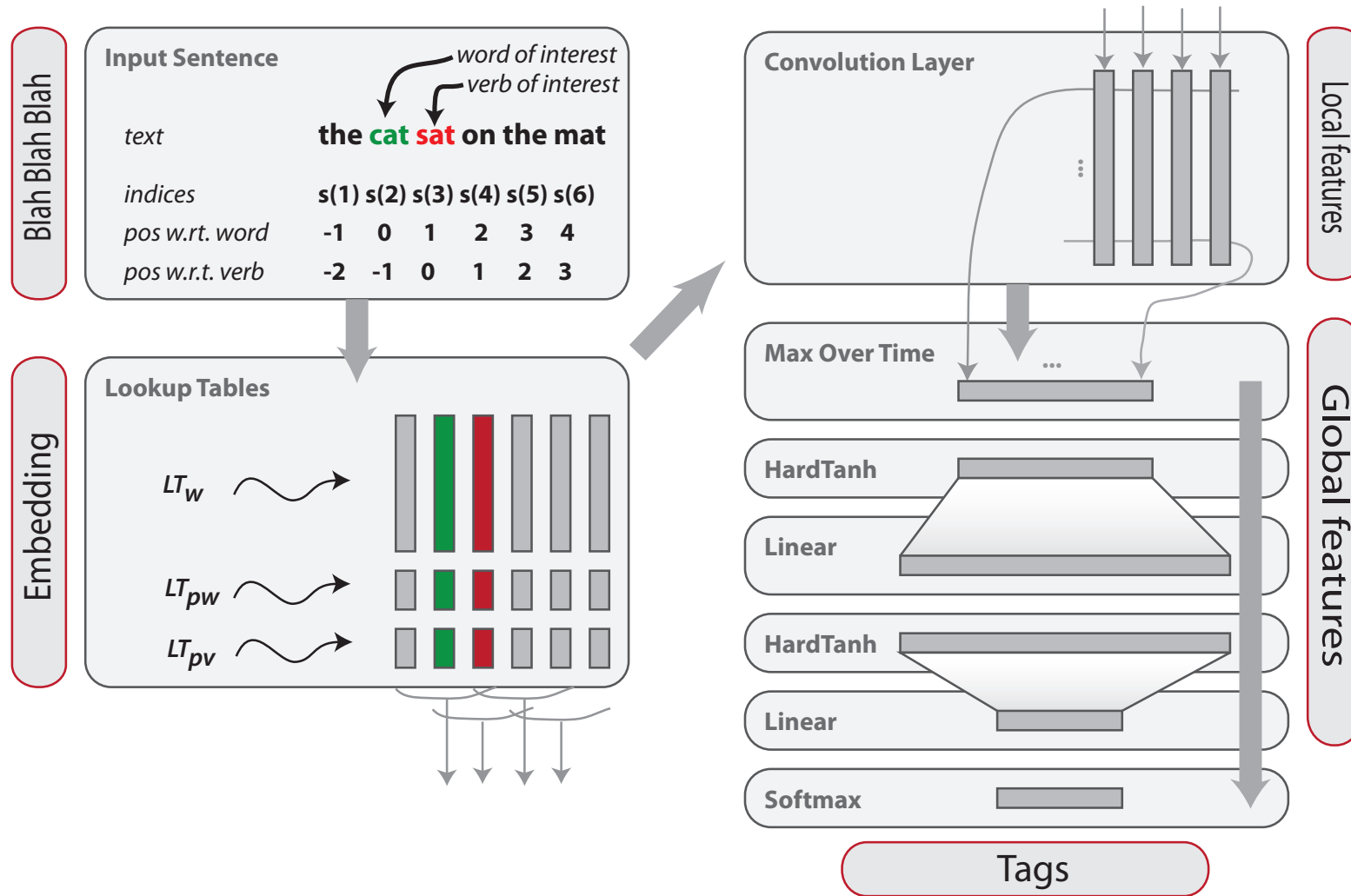
# The Deep Learning Way

(1/2)



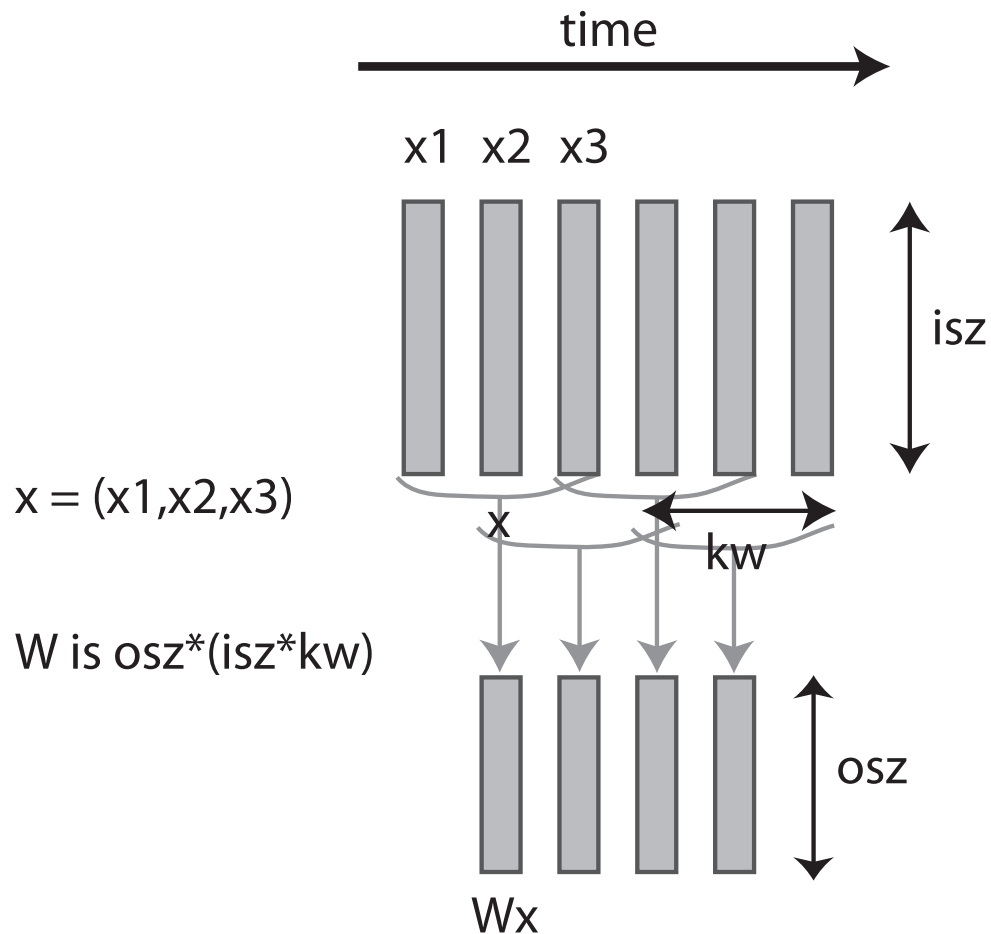
# The Deep Learning Way

(2/2)





# Convolutions



Extract **local** features – **share weights** through time/space



Used with **success** in **image** (Le Cun, **1989**) and **speech** (Bottou & Haffner, **1989**)

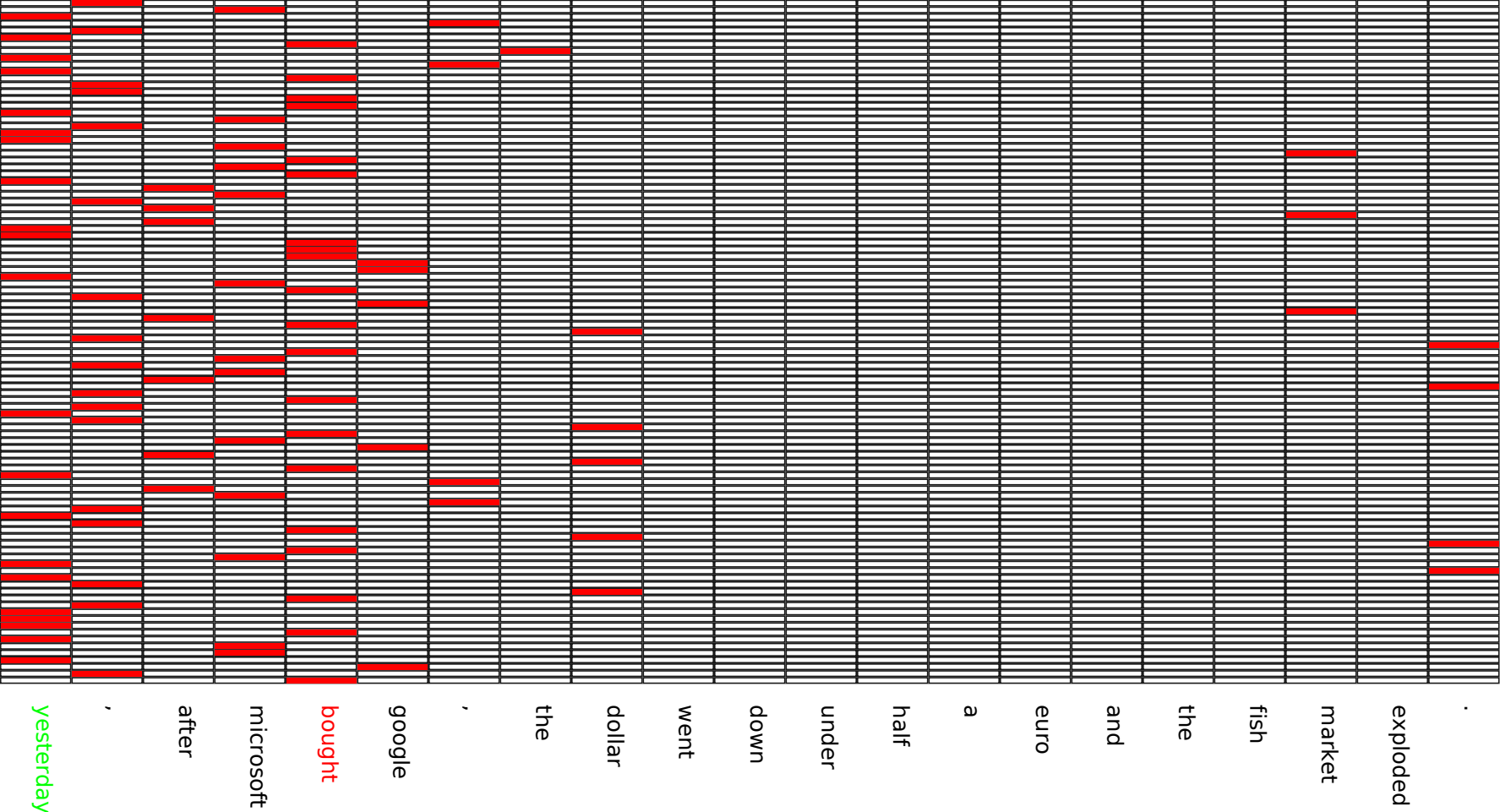


**Lookup-table** is a **special case**: convolution with kernel size of 1 and input  $i^{th}$  word

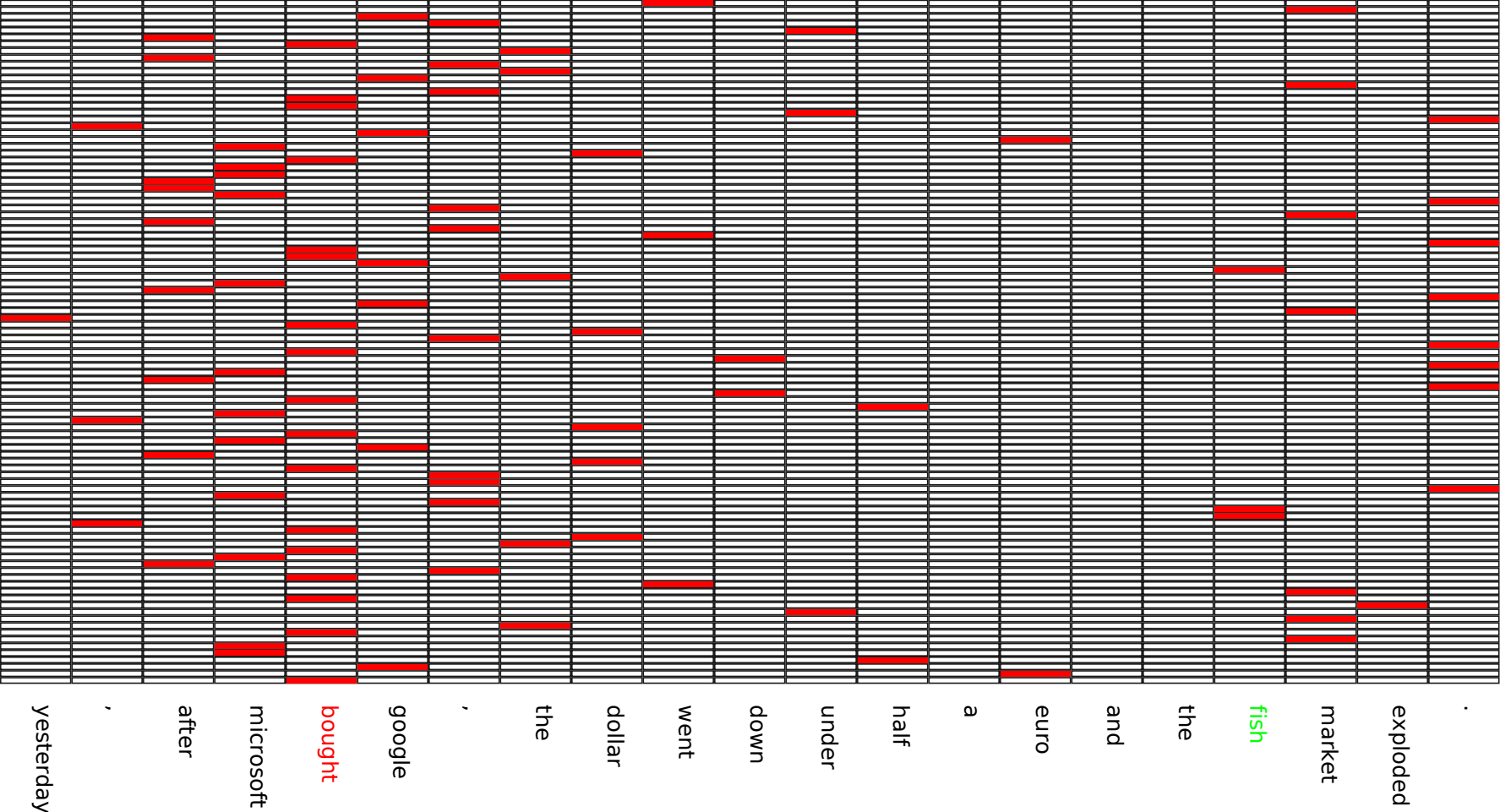
$(0, 0, \dots, 1, 0, \dots, 0)$  1 at position  $i$

Bengio et al (2001)

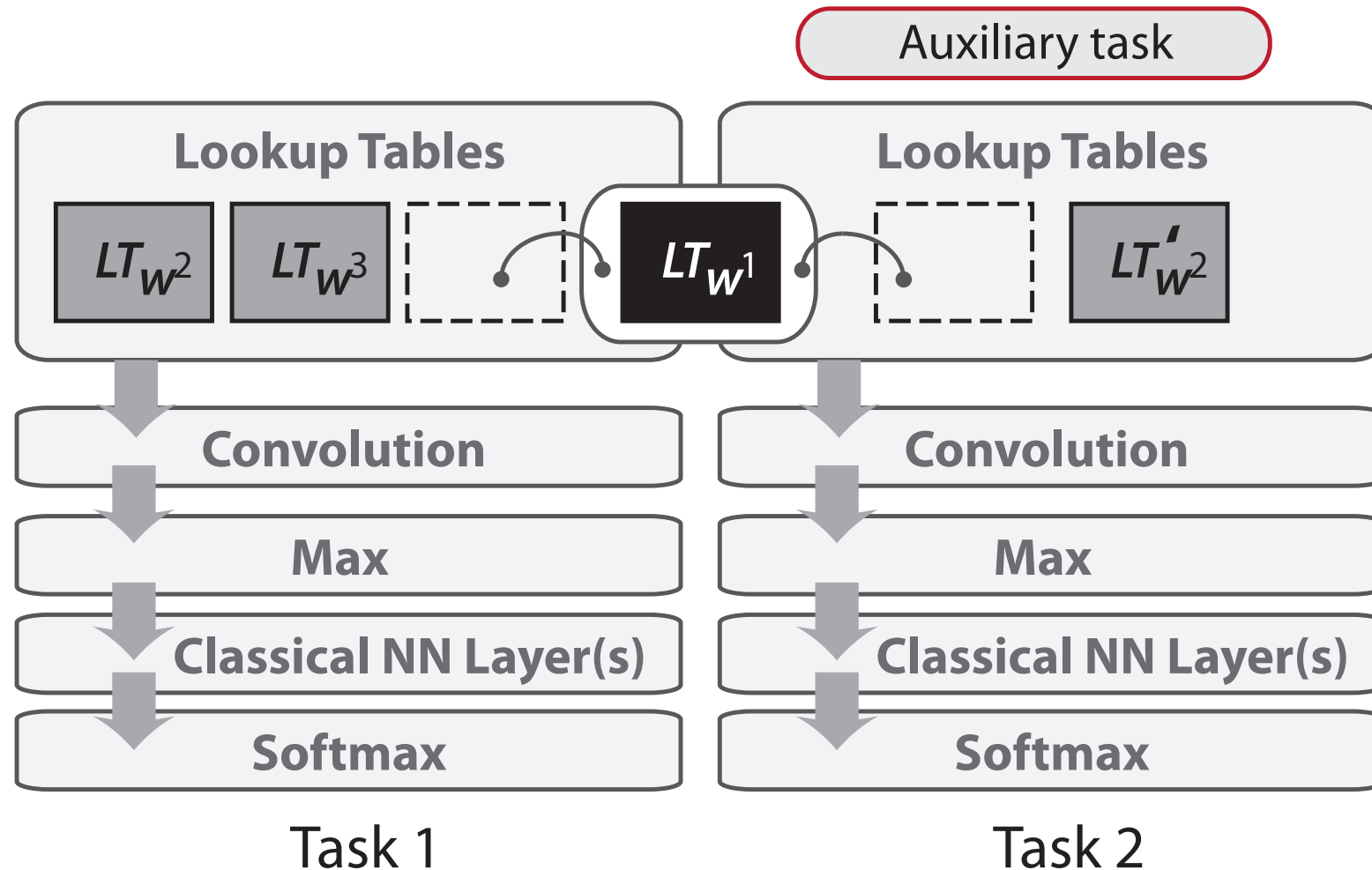
# Removing The Time Dimension (1/2)



# Removing The Time Dimension (2/2)



# Multi-Task Learning



Good overview in Caruana (1997)

# Improving Word Embedding

---

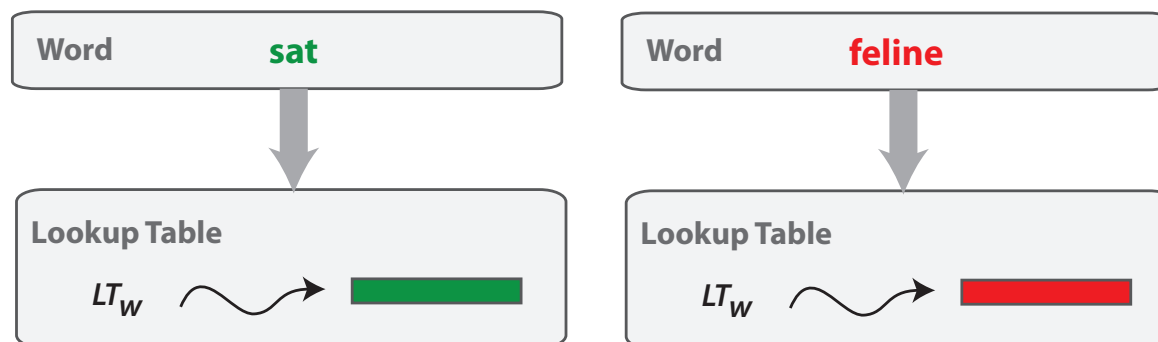


Rare words are **not trained** properly



Sentences with **similar words** should be **tagged in the same way**:

- ★ The **cat** sat on the mat
- ★ The **feline** sat on the mat



Wordnet

- ★ **pull** together linked words
- ★ **push** apart other pair of words

# Language Model: Think Massive

---



Language Model: “*is a sentence actually english or not?*”

Implicitly captures:   ★ syntax   ★ semantics



Bengio & Ducharme (2001) Probability of next word given previous words. Overcomplicated – we do not need probabilities here



English sentence windows: Wikipedia (~ 631M words)

Non-english sentence windows: middle word randomly replaced



Multi-class margin cost:

$$\sum_{s \in \mathcal{S}} \sum_{w \in \mathcal{D}} \max(0, 1 - f(s, w_s^*) + f(s, w))$$

$\mathcal{S}$ : sentence windows    $\mathcal{D}$ : dictionary

$w_s^*$ : true middle word in  $s$

$f(s, w)$ : network score for sentence  $s$  and middle word  $w$

# Language Model: Embedding

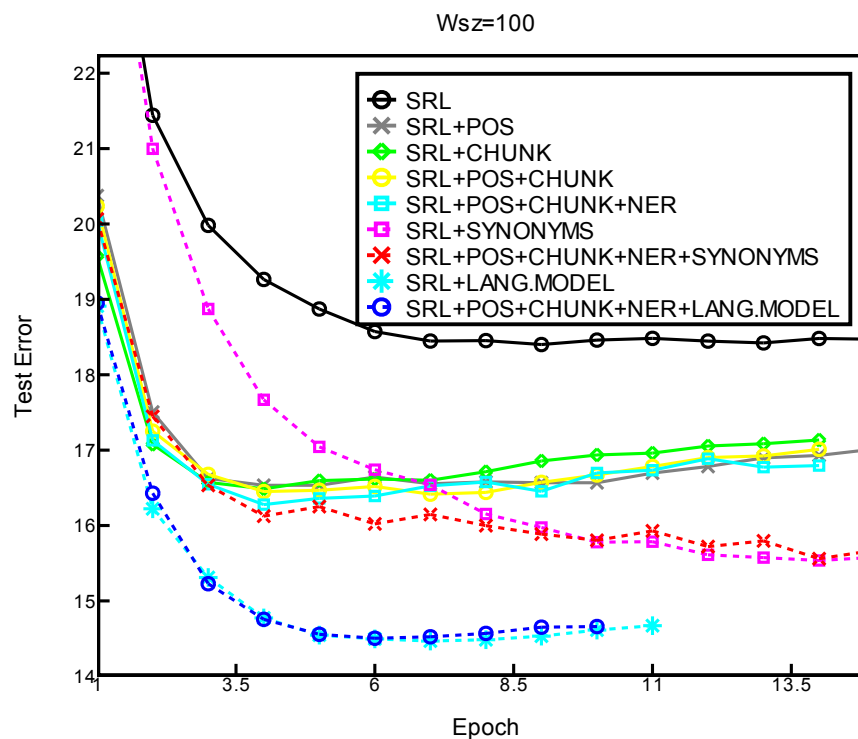
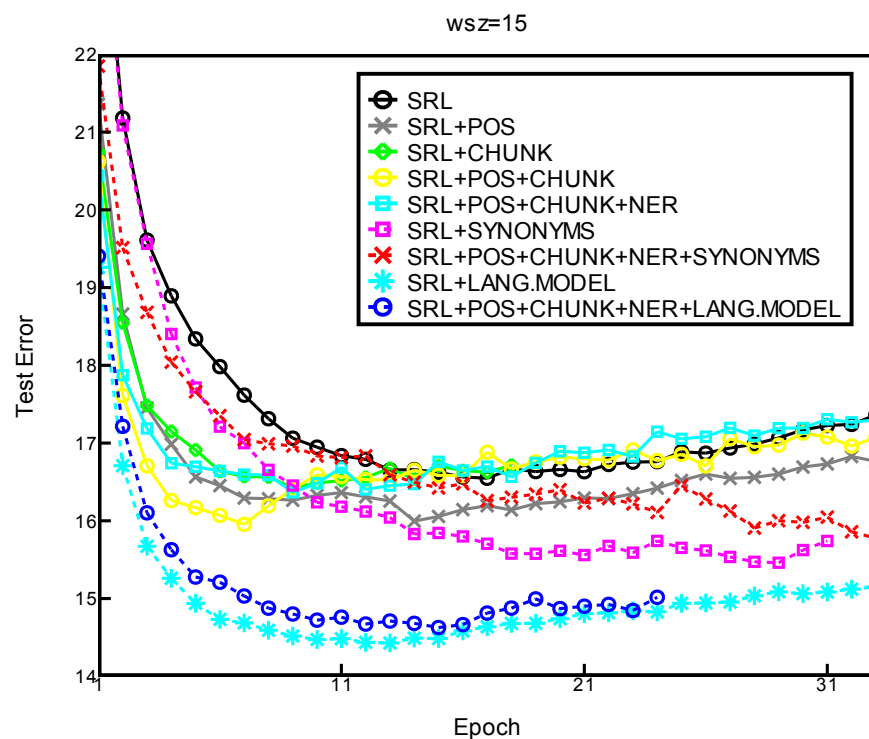
---

france 454	jesus 1973	xbox 6909	reddish 11724	scratched 29869
spain	christ	playstation	yellowish	smashed
italy	god	dreamcast	greenish	ripped
russia	resurrection	psNUMBER	brownish	brushed
poland	prayer	snest	bluish	hurled
england	yahweh	wii	creamy	grabbed
denmark	josephus	nes	whitish	tossed
germany	moses	nintendo	blackish	squeezed
portugal	sin	gamecube	silvery	blasted
sweden	heaven	psp	greyish	tangled
austria	salvation	amiga	paler	slashed

---

Dictionary size: 30,000 words. Even rare words are well embedded.

# MTL: Semantic Role Labeling



We get: 14.30%. State-of-the-art: 16.54% – Pradhan et al. (2004)



250× faster than state-of-the-art. ~ 0.01s to label a WSJ sentence.



# MTL: Unified Network for NLP

---

Improved results with Multi-Task Learning (MTL)

Task	Alone	MTL
SRL	18.40%	14.30%
POS	2.95%	2.91%
Chunking – error rate	5.4%	4.9%
Chunking – F1-score	91.5%	93.6%



POS: state-of-the-art  $\sim 3\%$



Chunking: Best system had 93.48% F1-score at CoNLL-2000 challenge <http://www.cnts.ua.ac.be/conll2000/chunking>. State-of-the-art is 94.1%. We get 94.9% by using POS features.

# Summary

---



We developed a **deep neural network architecture** for NLP



## Advantages

- ★ **General** to any NLP tagging task
- ★ **State-of-the-art** performance
- ★ **No hand designed features**
- ★ **Joint training**
- ★ Can exploit **massive unlabeled data**
- ★ **Extremely fast**: 0.02s for all tags of a sentence



## Inconvenients

- ★ Neural networks are a **powerful** tool: **hard to handle**



## Early Impacts

- ★ Easy to apply to other tasks or languages: extending to **Japanese**
- ★ Fast: developed a **semantic search** system