

Deep Learning for Text

From Word Embeddings to Convolutional Neural Networks

Martin Jaggi

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



SPINNINGBYTES

SwissText Conference, 8th June 2016

Natural Language Processing



- ❖ Numerous applications with huge impact:
 - ❖ Search - access to information
 - ❖ Question answering - access to knowledge
 - ❖ Machine translation - bridge multi-linguality
 - ❖ Machine reading & summarization - essence of text
 - ❖ Conversational agents - talk the talk
- ❖ ... we are only at the beginning!

Outline

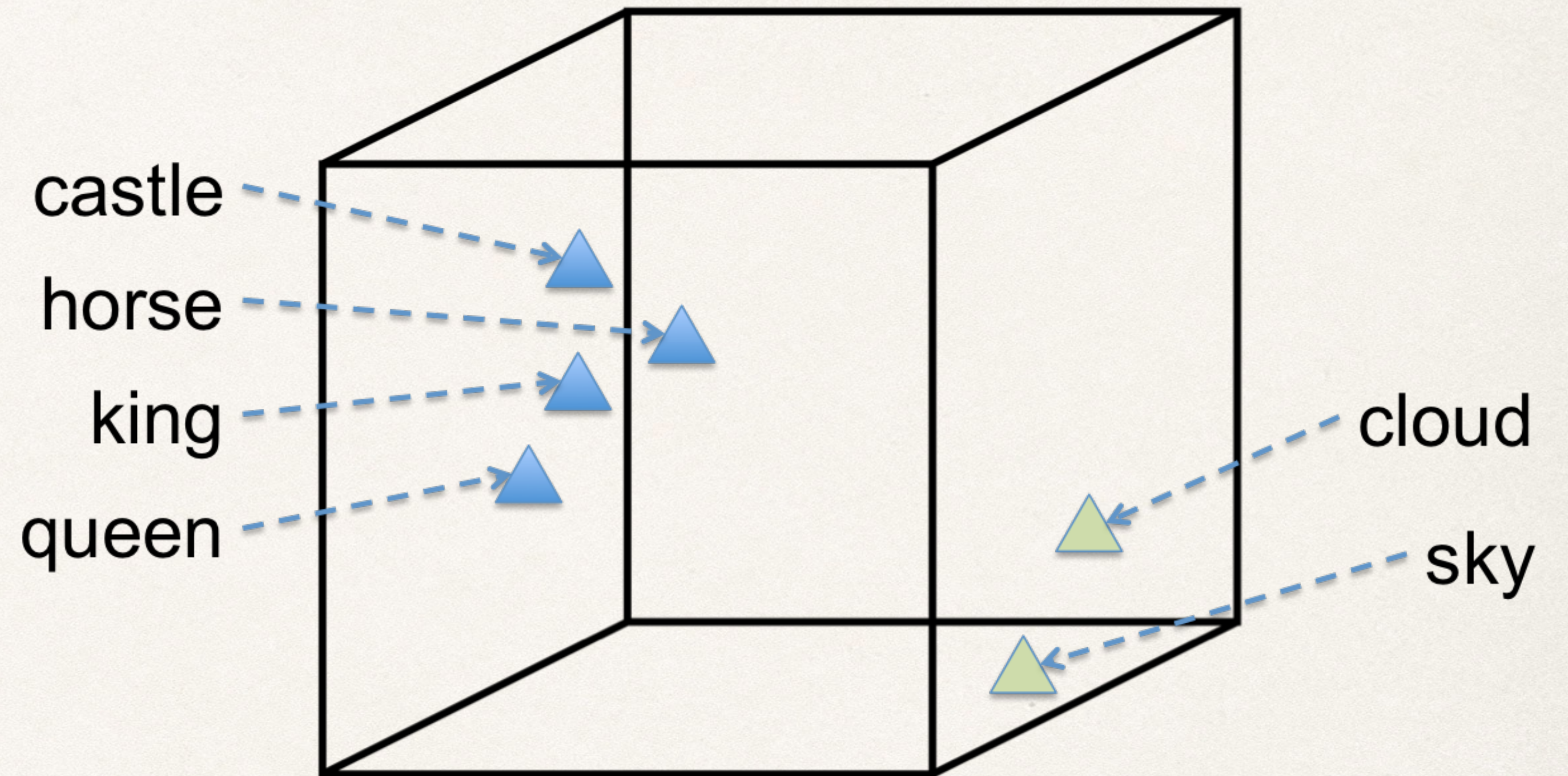
- ❖ Semantic Text Representations
 - ❖ Word Embeddings
 - ❖ Document Embeddings
 - ❖ Applications of Machine Learning to Text

From Words to Features

- ✦ Bag of words representation

$$i \longrightarrow (0, \dots, 1, \dots, 0) \in \mathbb{R}^{1M}$$

Word Embeddings



$$i \longrightarrow \mathbf{v}_i \in \mathbb{R}^{50}$$

Word Embeddings

Word

	1	1		
		3		
	1			
	2		1	
1				1
		1		
	1		1	1

explain co-occurrence i, j
by means of

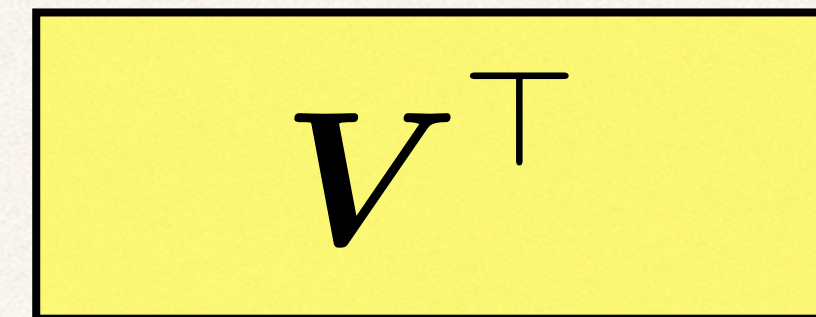
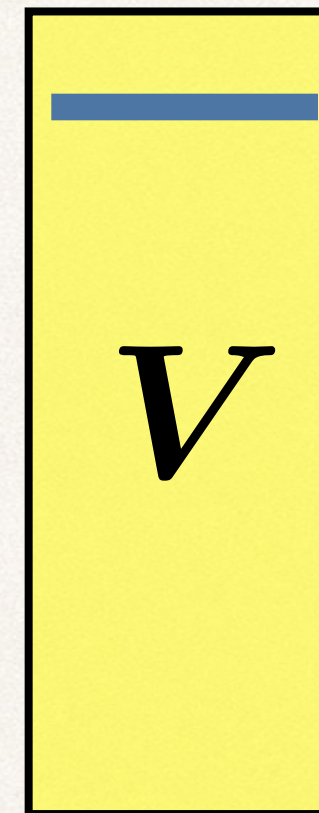
$$\mathbf{v}_i^\top \mathbf{v}_j$$

Word Embeddings

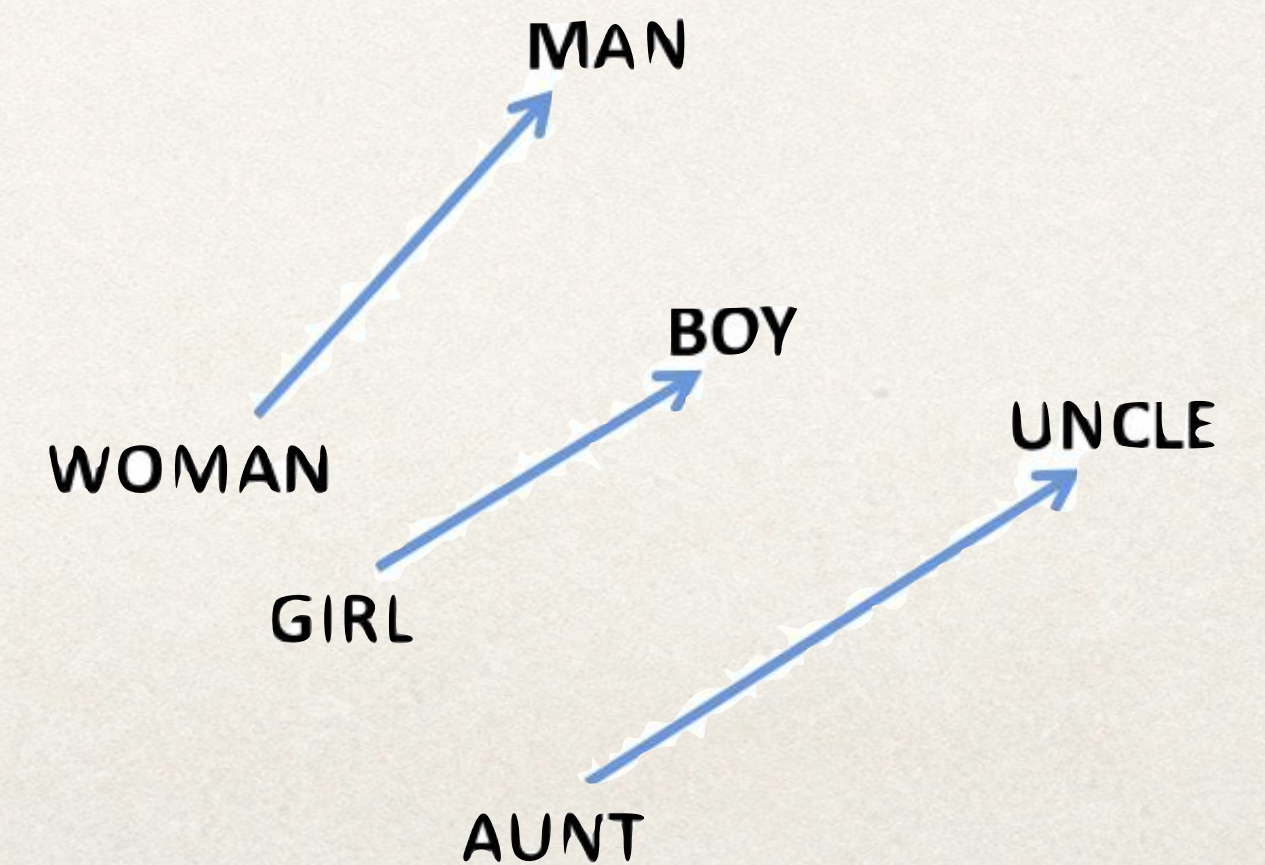
Word

	1	1		
		3		
	1			
	2		1	
1				1
		1		
	1		1	1

\approx



SVD, PLSA etc
word2vec, gloVe



Word Embeddings



Movies

Customers

	★	★ ★ ★		
		★ ★ ★ ★		
	★			
	★ ★		★ ★ ★	
★ ★ ★ ★				★ ★ ★
		★ ★		
	★ ★		★	★ ★ ★

$$\approx UV^T$$

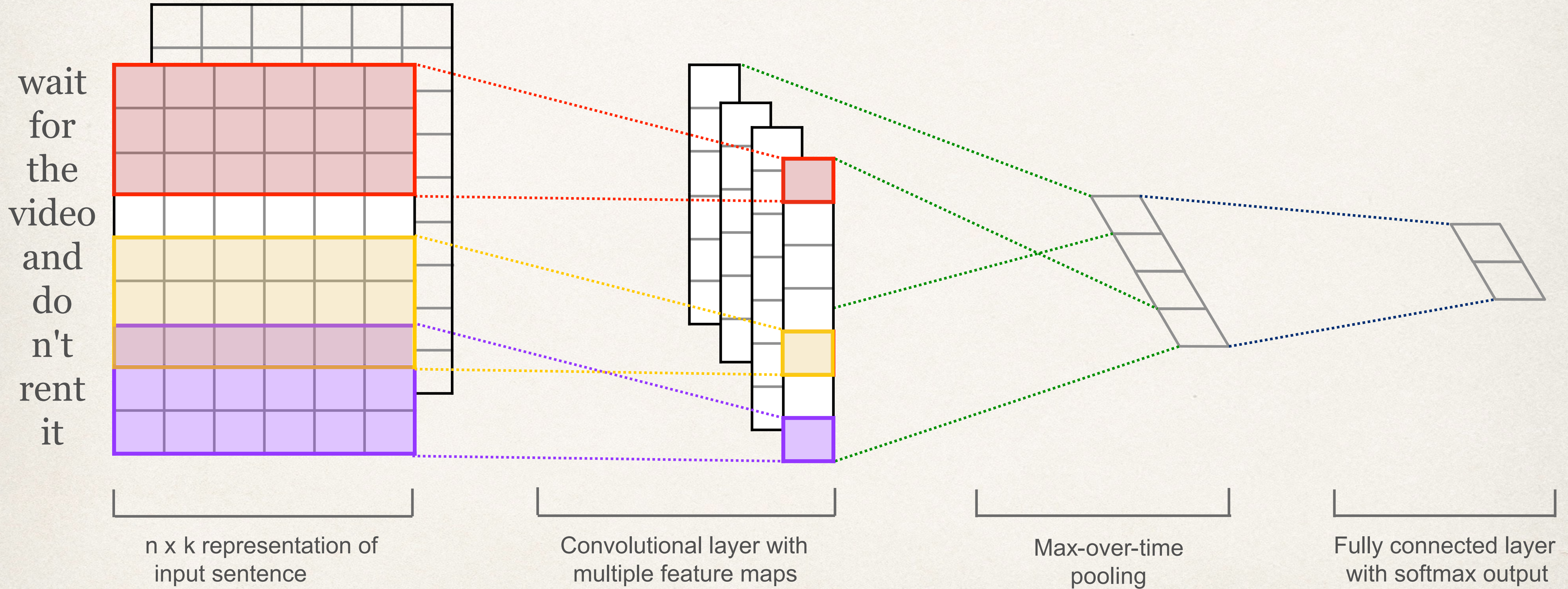
Word Embeddings - Summary

- ❖ Very successful new variation of an old theme
- ❖ State of the art feature representations for words
- ❖ Not related to deep learning
- ❖ Parallelization still challenging
- ❖ Limited to represent words or short n -grams

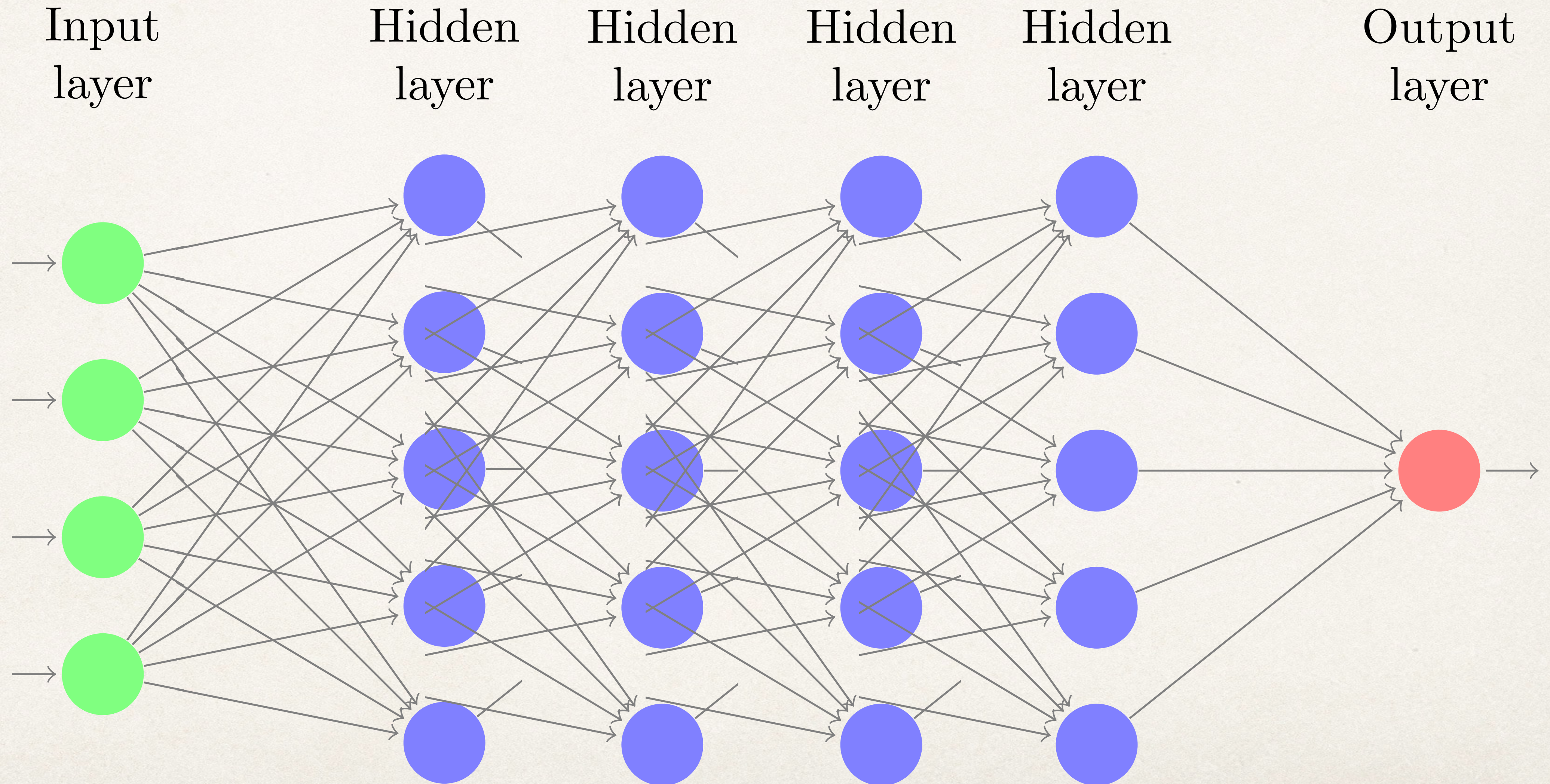
Document Embeddings

- ❖ How to represent a sequence of words?

A Two-Stage Approach



Neural Networks



Alternative Document Embeddings

- ❖ Convolutional Neural Networks (CNN)
- ❖ Long Short-Term Memory (LSTM) Networks
- ❖ paragraph2vec / doc2vec

Application: Sentiment Classification

- ❖ A state-of-the art system for text classification
- ❖ Two ETH Master Theses by

Jan Deriu & Maurice Gonzenbach

Results

- ❖ **SemEval Competition**
 - ❖ running since 1998
 - ❖ new set of manually annotated tweets every year
- ❖ **Our Entries in the Sentiment Competition**
 - ❖ 2016 **1st place** (Convolutional NN, ensemble)
 - ❖ 2015 8th place (SVM, lexica, ensemble)
 - ❖ 2014 8th place (SVM, lexica, ensemble)

Results

- ❖ SemEval Competition

 - ❖ running since 1998

 - ❖ new set of manually annotated

- ❖ Our Entries in the Sentiment Comp

 - ❖ 2016 1st place (Convolutional N

 - ❖ 2015 8th place (SVM, lexica, ens

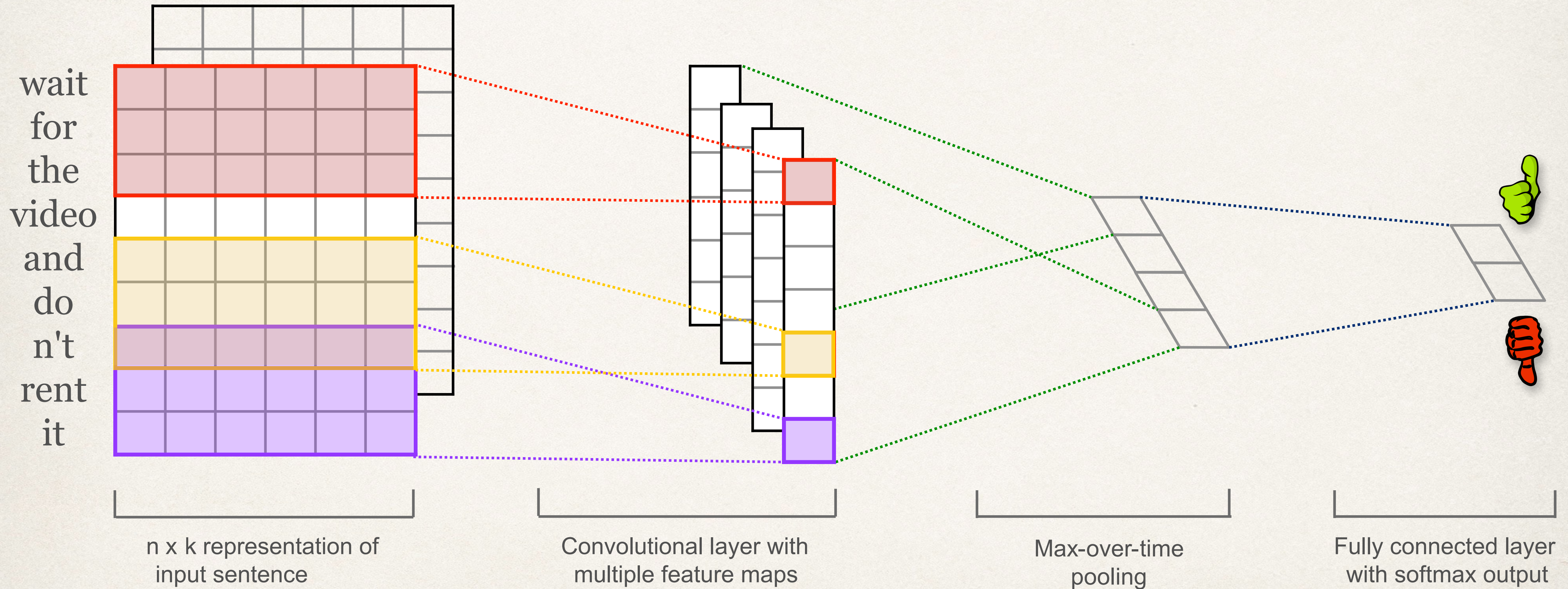
 - ❖ 2014 8th place (SVM, lexica, ens



<i>negative</i>	<i>neutral</i>	But i wanna wear my Concords tomorrow though but i don't feel like it
<i>positive</i>	<i>neutral</i>	Gonna watch Grey's Anatomy all day today and tomorrow(:
<i>negative</i>	<i>neutral</i>	@CoachVac heey do you know anything about UVA's fallll fest loll they
<i>neutral</i>	<i>neutral</i>	@DustyEf when that sun is high in that Texas sky, I'll be buckin it to co
<i>neutral</i>	<i>positive</i>	Up 20 points in my money league with Vernon Davis and L. Fitz still to
<i>neutral</i>	<i>positive</i>	DEEJAYING this FRIDAY in THE FIRST CHOP it's CHRIS actual SMITH w
<i>negative</i>	<i>negative</i>	The Rick Santorum signing that was scheduled for tomorrow at the Boc
<i>positive</i>	<i>neutral</i>	@dreami9 lol yep looks like it! Was after El Clasico on Sunday. I didn't
<i>neutral</i>	<i>neutral</i>	Back in Stoke on Trent for the 2nd time today!
<i>neutral</i>	<i>neutral</i>	First Girls Varsity Basketball Game tomorrow at 6:00 pm Then Football
<i>neutral</i>	<i>neutral</i>	#UFC lightweights @Young__Assassin VS @jamievarner set for TUF 16
<i>neutral</i>	<i>neutral</i>	@OOOOO_WEEEE slide thru sometime this weekend ill have somethin v
<i>negative</i>	<i>negative</i>	@DannyB618 Sure absolutely-- I meant out of the Bachmann, Perry, S
<i>negative</i>	<i>negative</i>	@RichardGordon48 re Levein discussion on Wed. Can't keep changing b
<i>neutral</i>	<i>neutral</i>	Today In History November 02, 1958 Elvis gave a party at his hotel bef
<i>neutral</i>	<i>positive</i>	Hustle cause you got to then kick back n party everyday like its Fri
<i>positive</i>	<i>positive</i>	I can't sleep. Way too exited about Vancouver tomorrow! I'm like a kid

#	System	2013		2014			2015	2016
		Tweet	SMS	Tweet	Tweet sarcasm	Live- Journal	Tweet	Tweet
1	SwissCheese	0.700 ₄	0.637 ₂	0.716 ₄	0.566 ₁	0.695 ₇	0.671 ₁	0.633 ₁
2	SENSEI-LIF	0.706 ₃	0.634 ₃	0.744 ₁	0.467 ₈	0.741 ₁	0.662 ₂	0.630 ₂
3	UNIMELB	0.687 ₆	0.593 ₉	0.706 ₆	0.449 ₁₁	0.683 ₉	0.651 ₄	0.617 ₃
4	INESC-ID	0.723 ₁	0.609 ₆	0.727 ₂	0.554 ₂	0.702 ₄	0.657 ₃	0.610 ₄
5	aueb.twitter.sentiment	0.666 ₇	0.618 ₅	0.708 ₅	0.410 ₁₇	0.695 ₇	0.623 ₇	0.605 ₅
6	SentiSys	0.714 ₂	0.633 ₄	0.723 ₃	0.515 ₄	0.726 ₂	0.644 ₅	0.598 ₆
7	I2RNTU	0.693 ₅	0.597 ₇	0.680 ₇	0.469 ₆	0.696 ₆	0.638 ₆	0.596 ₇
8	INSIGHT-1	0.602 ₁₆	0.582 ₁₂	0.644 ₁₅	0.391 ₂₃	0.559 ₂₃	0.595 ₁₆	0.593 ₈
9	TwISE	0.610 ₁₅	0.540 ₁₆	0.645 ₁₃	0.450 ₁₀	0.649 ₁₃	0.621 ₈	0.586 ₉
10	ECNU (*)	0.643 ₉	0.593 ₉	0.662 ₈	0.425 ₁₄	0.663 ₁₀	0.606 ₁₁	0.585 ₁₀
11	NTNUSentEval	0.623 ₁₁	0.641 ₁	0.651 ₁₀	0.427 ₁₃	0.719 ₃	0.599 ₁₃	0.583 ₁₁
12	MDSSENT	0.589 ₁₉	0.509 ₂₀	0.587 ₂₀	0.386 ₂₄	0.606 ₁₈	0.593 ₁₇	0.580 ₁₂
	CUFE	0.642 ₁₀	0.596 ₈	0.662 ₈	0.466 ₉	0.697 ₅	0.598 ₁₄	0.580 ₁₂
14	THUIR	0.616 ₁₂	0.575 ₁₄	0.648 ₁₁	0.399 ₂₀	0.640 ₁₅	0.617 ₁₀	0.576 ₁₄
	PUT	0.565 ₂₁	0.511 ₁₉	0.614 ₁₉	0.360 ₂₇	0.648 ₁₄	0.597 ₁₅	0.576 ₁₄
16	LYS	0.650 ₈	0.579 ₁₃	0.647 ₁₂	0.407 ₁₈	0.655 ₁₁	0.603 ₁₂	0.575 ₁₆
17	IIP	0.598 ₁₇	0.465 ₂₃	0.645 ₁₃	0.405 ₁₉	0.640 ₁₅	0.619 ₉	0.574 ₁₇
18	UniPI	0.592 ₁₈	0.585 ₁₁	0.627 ₁₇	0.381 ₂₅	0.654 ₁₂	0.586 ₁₈	0.571 ₁₈
19	DIEGOLab16 (*)	0.611 ₁₄	0.506 ₂₁	0.618 ₁₈	0.497 ₅	0.594 ₂₀	0.584 ₁₉	0.554 ₁₉
20	GTI	0.612 ₁₃	0.524 ₁₇	0.639 ₁₆	0.468 ₇	0.623 ₁₇	0.584 ₁₉	0.539 ₂₀
21	OPAL	0.567 ₂₀	0.562 ₁₅	0.556 ₂₃	0.395 ₂₁	0.593 ₂₁	0.531 ₂₁	0.505 ₂₁
22	DSIC-ELIRF	0.494 ₂₅	0.404 ₂₆	0.546 ₂₆	0.342 ₂₉	0.517 ₂₄	0.531 ₂₁	0.502 ₂₂
23	UofL	0.490 ₂₆	0.443 ₂₄	0.547 ₂₅	0.372 ₂₆	0.574 ₂₂	0.502 ₂₅	0.499 ₂₃
	ELiRF	0.462 ₂₈	0.408 ₂₅	0.514 ₂₈	0.310 ₃₃	0.493 ₂₅	0.493 ₂₆	0.499 ₂₃
25	ISTI-CNR	0.538 ₂₂	0.492 ₂₂	0.572 ₂₁	0.327 ₃₀	0.598 ₁₉	0.508 ₂₄	0.494 ₂₅
26	SteM	0.518 ₂₃	0.315 ₂₉	0.571 ₂₂	0.320 ₃₂	0.405 ₂₈	0.517 ₂₃	0.478 ₂₆
27	Tweester	0.506 ₂₄	0.340 ₂₈	0.529 ₂₇	0.540 ₃	0.379 ₂₉	0.479 ₂₈	0.455 ₂₇
28	Minions	0.489 ₂₇	0.521 ₁₈	0.554 ₂₄	0.420 ₁₆	0.475 ₂₆	0.481 ₂₇	0.415 ₂₈
29	Aicyber	0.418 ₂₉	0.361 ₂₇	0.457 ₂₉	0.326 ₃₁	0.440 ₂₇	0.432 ₂₉	0.402 ₂₉
30	mib	0.394 ₃₀	0.310 ₃₀	0.415 ₃₁	0.352 ₂₈	0.359 ₃₁	0.413 ₃₁	0.401 ₃₀
31	VCU-TSA	0.383 ₃₁	0.307 ₃₁	0.444 ₃₀	0.425 ₁₄	0.336 ₃₂	0.416 ₃₀	0.372 ₃₁
32	SentimentalITists	0.339 ₃₃	0.238 ₃₃	0.393 ₃₂	0.288 ₃₄	0.323 ₃₄	0.343 ₃₃	0.339 ₃₂
33	WR	0.355 ₃₂	0.284 ₃₂	0.393 ₃₂	0.430 ₁₂	0.366 ₃₀	0.377 ₃₂	0.330 ₃₃
34	CICBUAPnlp	0.193 ₃₄	0.193 ₃₄	0.335 ₃₄	0.393 ₂₂	0.326 ₃₃	0.303 ₃₄	0.303 ₃₄

Convolutional Neural Network (CNN)

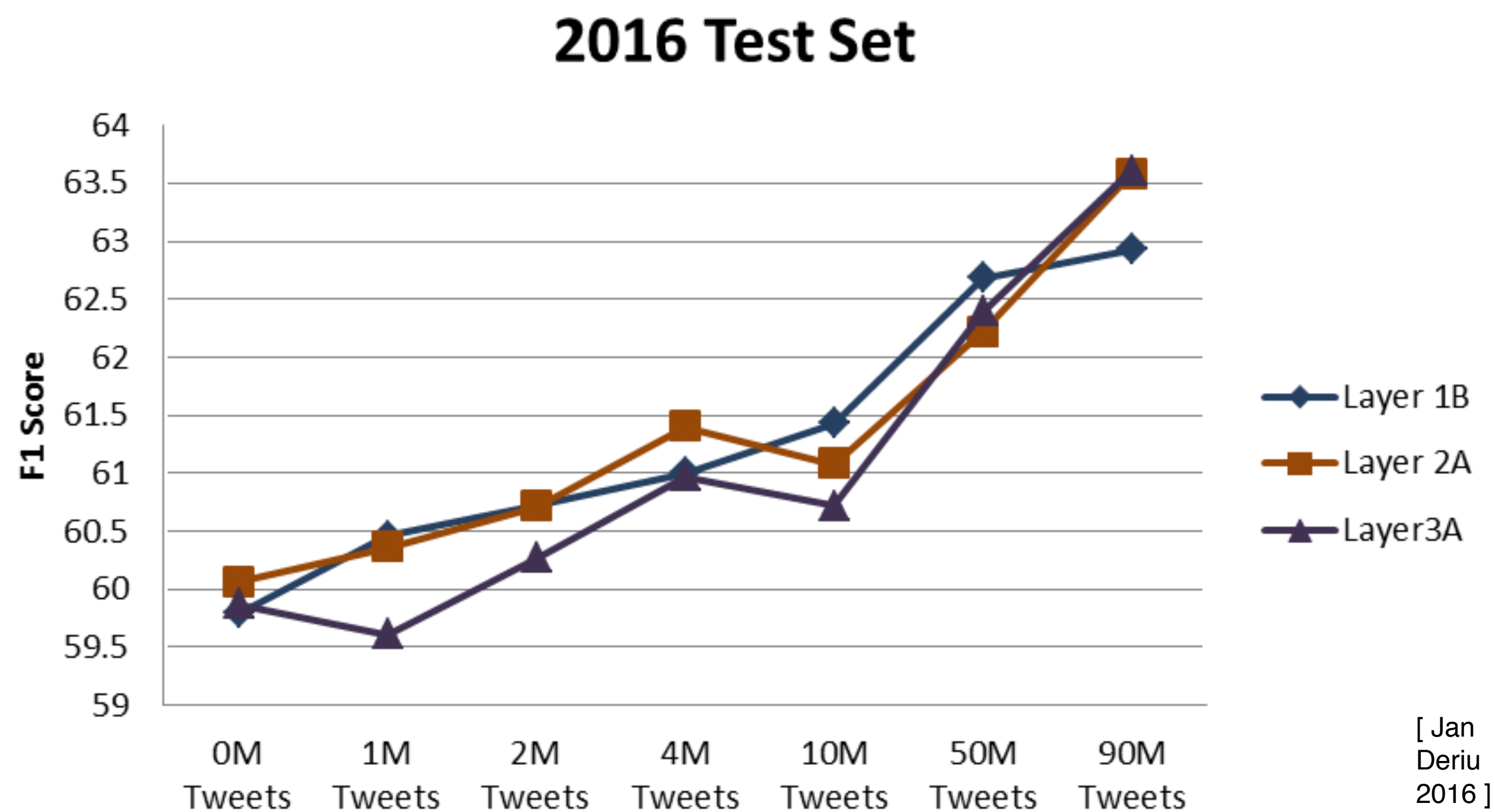


Distant Supervision

- ❖ 90M tweets containing :) or :(



System based on
(Severyn & Moschitti, 2015)



Evolving Word Embeddings

- ❖ Backpropagation changes both NN weights & word embeddings



Outlook

- ❖ Sentence / document embeddings are useful for many tasks
- ❖ Medical applications
- ❖ Depth of the NNs?
- ❖ **Un**-supervised training?

References

- ❖ Many online resources, open source frameworks etc, active community
- ❖ Master Theses Jan Deriu & Maurice Gonzenbach
- ❖ *SwissCheese at SemEval-2016 Task 4: Sentiment Classification Using an Ensemble of Convolutional Neural Networks with Distant Supervision*

Thanks

Jan Deriu, Maurice Gonzenbach, Fatih Uzdilli, Aurelien Lucchi, Valeria De Luca, Dominic Egger, Pascal Julmy, Leon Derczynski, Mark Cieliebak