

# $\ell_1$ Regularization of Word Embeddings for Multi-Word Expression Identification

Gábor Berend

Institute of Informatics, University of Szeged  
berendg@inf.u-szeged.hu



## Overview

- Multi-word expressions (MWEs) are lexical items with spaces (called *catenas* in linguistics)
- We conducted controlled experiments to compare the effects of various **word representations** and **classification algorithms** for MWE identification

## MWE dataset used

- Wiki150 corpus Vincze et al. (2011) (cca. 114K tokens and 4.4K sentences)
- Multiple types of MWEs

MWE type	Example
Noun compounds	<i>black box</i>
Adjectival compounds	<i>monkey styled</i>
Verb-Particle Constructions (VPC)	<i>went on</i>
Light-Verb Constructions (LVC)	<i>opens fire</i>
Idioms	<i>caught the eye of</i>
Other	<i>alter ego</i>
Location NE	<i>Sierra Leone</i>
Person NE	<i>Sir Elton John</i>
Organization NE	<i>Major Indoor Soccer League</i>

## Compared models

- Linear CRF model based on pre-trained dense word embeddings
  - Glove, polyglot, skip-gram, CBOW embeddings trained on English Wikipedia
- Linear CRF model based on  $\ell_1$ -regularized word embeddings
  - For an embedding matrix  $X \in \mathbb{R}^{d \times |V|}$  we solve for  $\min_{D \in \mathcal{C}, \alpha} \|X - D\alpha\|_F^2 + \lambda \|\alpha\|_1$ 
    - $\mathcal{C}$  is the convex set of  $\mathbb{R}^{d \times k}$  matrices with column norms  $\leq 1$ , and  $\alpha$  contains the sparse coefficients for word forms
    - For word  $i$  features are derived by  $\phi_{sparse}(i) = \{j : 1|\alpha_x[j, i] > 0 \wedge 1 \leq j \leq k\}$
- Linear CRF model based on Brown-cluster prefixes of word forms as features
- CRF model based on “traditional” feature set ( $\oplus$  denoting concatenation)

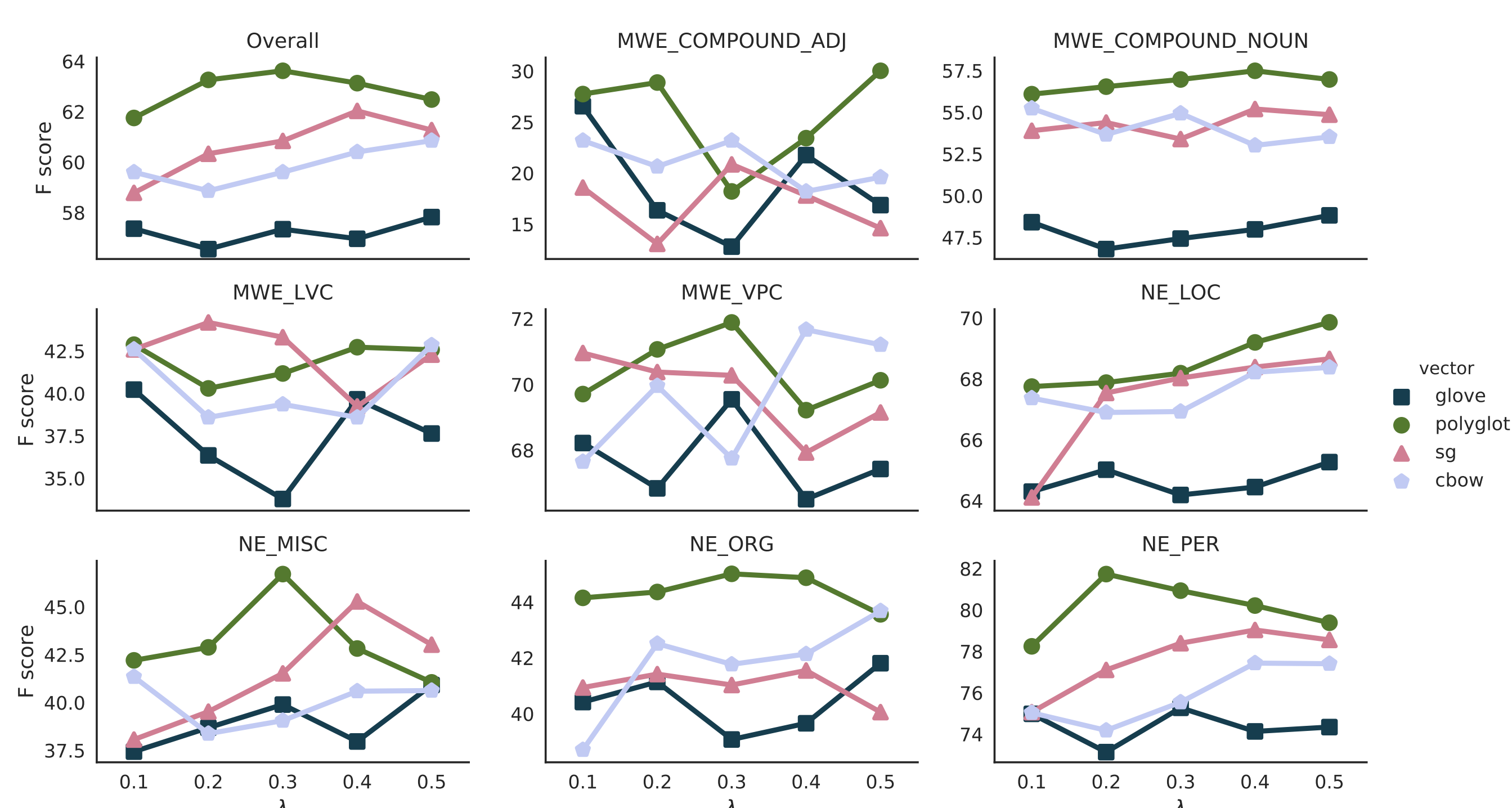
### Feature template

$\text{isNumber}(w_t)$	
$\text{isTitleCase}(w_t)$	
$\text{isNonAlnum}(w_t)$	
$\text{prefix}(w_t, i)$	$1 \leq i \leq 4$
$\text{suffix}(w_t, i)$	$1 \leq i \leq 4$
$w_{t+j}$	$-2 \leq j \leq 2$
$w_t \oplus w_{t+j}$	$1 \leq j \leq 9$
$w_t \oplus w_{t-j}$	$1 \leq j \leq 9$
$\bigoplus_{i=t+j}^{t+j+1} w_i$	$-2 \leq j \leq 1$
$\bigoplus_{i=t+j}^{t+j+2} w_i$	$-2 \leq j \leq 0$
$\bigoplus_{i=t+j-1}^{t+j+2} w_i$	$-1 \leq j \leq 0$
$\bigoplus_{i=t-2}^{t+2} w_i$	

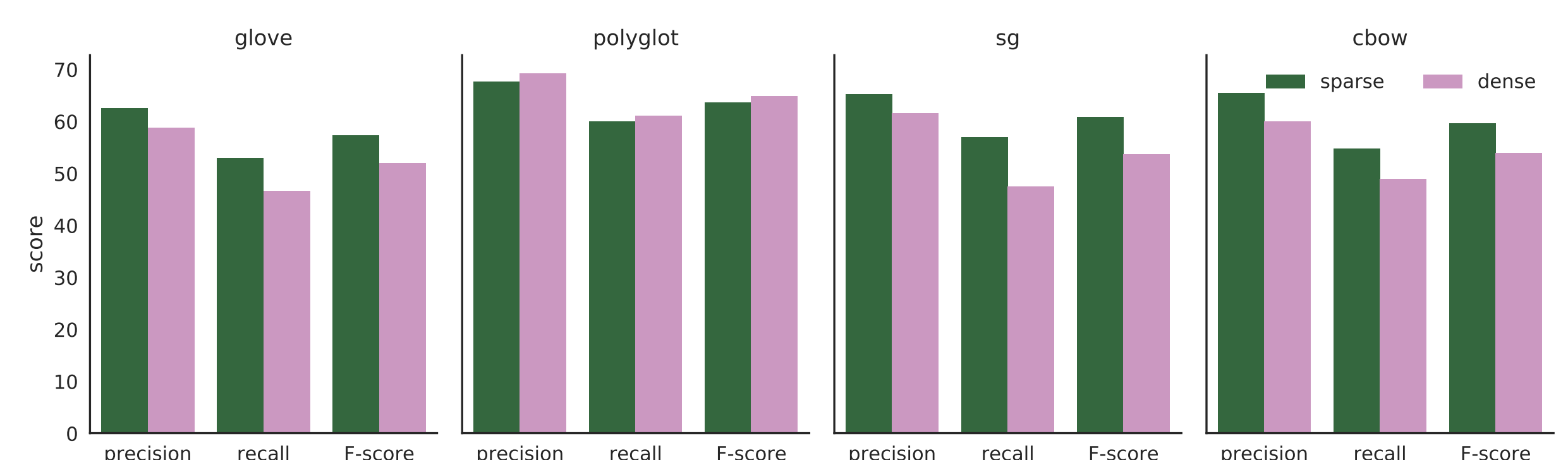
- word-level biLSTM models (optionally extended with character-level representations) (Plank et al., 2016)
  - Word embeddings initialized with polyglot (important due to the size of Wiki150)

## The effects of $\ell_1$ regularization

- Compared the effects of relying on differently trained word embeddings (all trained on the same Wikipedia dumps)
- $d = 64, k = 1000$ , sparsity ranging between 0.5% and 5%



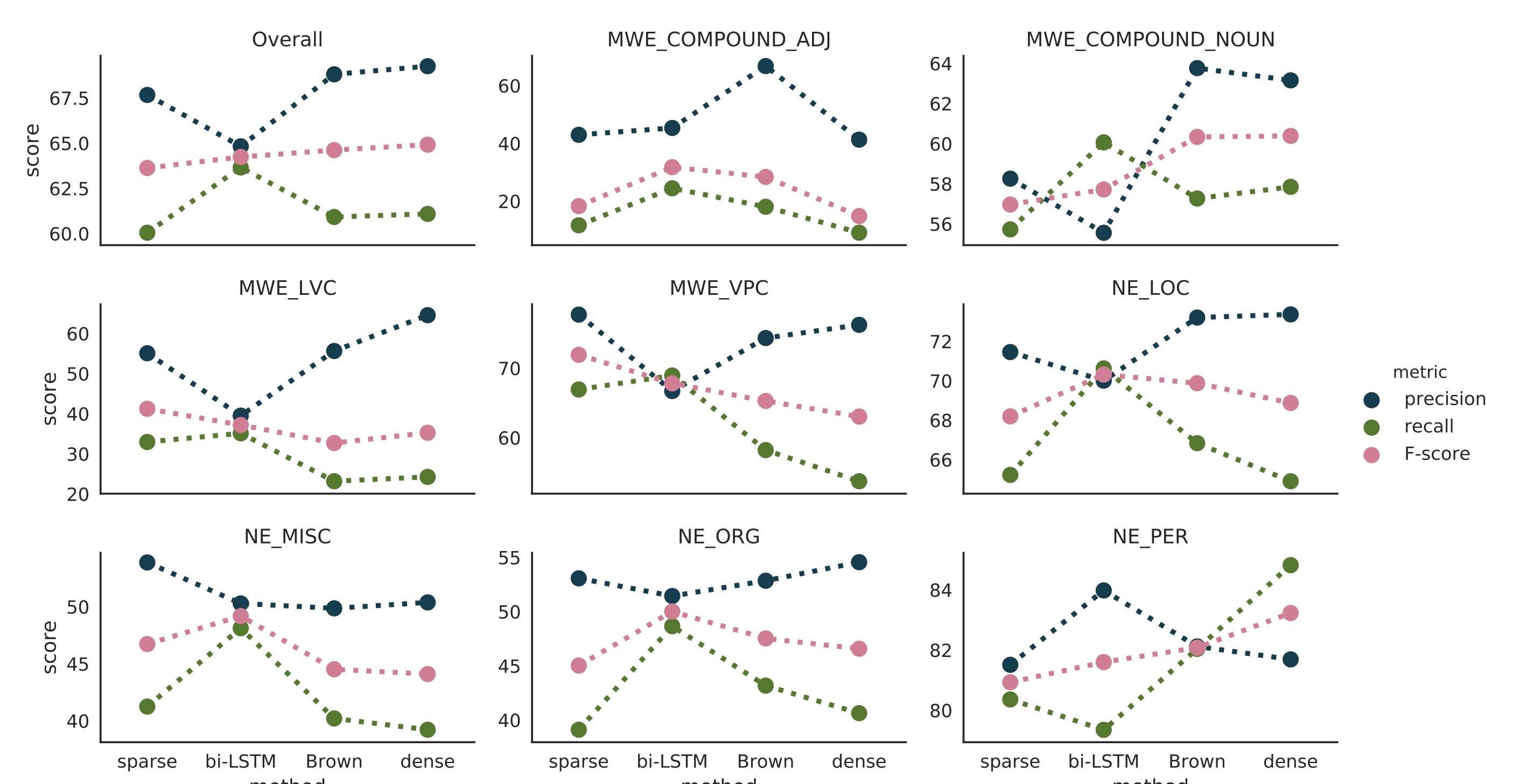
## Dense versus sparse representation



## Comparison of alternative models

method	Precision	Recall	F-score
FRw	62.91	18.04	28.04
FRwc	62.92	50.30	55.91
polyglot sparse ( $\lambda = 0.3$ )	67.65	60.03	63.61
bi-LSTM	64.81	63.64	64.22
Brown	68.79	60.90	64.60
polyglot dense	69.24	61.07	64.90

## Detailed performance of alternative models



## Extending biLSTM with character-level representation

- Improvements mostly on the NE categories

	biLSTM	biLSTM with char
Compound adj	<b>31.67</b>	6.32
Compound noun	57.70	<b>59.82</b>
LVC	<b>37.12</b>	33.74
VPC	<b>67.77</b>	61.17
Location	70.30	<b>72.98</b>
Misc	49.15	<b>51.66</b>
Organization	49.98	<b>53.55</b>
Person	81.59	<b>84.50</b>
Avg.	64.22	<b>66.48</b>

## Conclusions

- biLSTM tends to produce more balanced results with respect to precision-recall (except for the Persons)
- Models employing sparse word representation perform best on verb-related MWE types (VPC and LVC)



## References

- G. Berend. Sparse coding of neural word embeddings for multilingual sequence labeling. *Transactions of the Association for Computational Linguistics*, 5:247–261, 2017. ISSN 2307-387X. URL <https://transacl.org/ojs/index.php/tacl/article/view/1063>.
- B. Plank, A. Søgaard, and Y. Goldberg. Multilingual part-of-speech tagging with bidirectional long short-term memory models and auxiliary loss. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 412–418. Association for Computational Linguistics, 2016. URL <http://anthology.aclweb.org/P16-2067>.
- V. Vincze, I. Nagy T., and G. Berend. Multiword expressions and named entities in the Wiki50 corpus. In *Proceedings of the International Conference Recent Advances in Natural Language Processing 2011*, pages 289–295. Hissar, Bulgaria, September 2011. RANLP 2011 Organising Committee. URL <http://aclweb.org/anthology/R11-1040>.