

# Gating Mechanisms for Combining Character and Word-level Word Representations: An Empirical Study

Jorge A. Balazs, Yutaka Matsuo

The University of Tokyo, Graduate School of Engineering, Japan

{jorge, matsuo}@weblab.t.u-tokyo.ac.jp



## Problem

Incorporating subword information into word representations has been shown to be beneficial, however there is no principled way for doing so.

## Questions

- How does the method for combining character and word representations affect the quality of final word representations?
- What is effect these have in downstream performance?

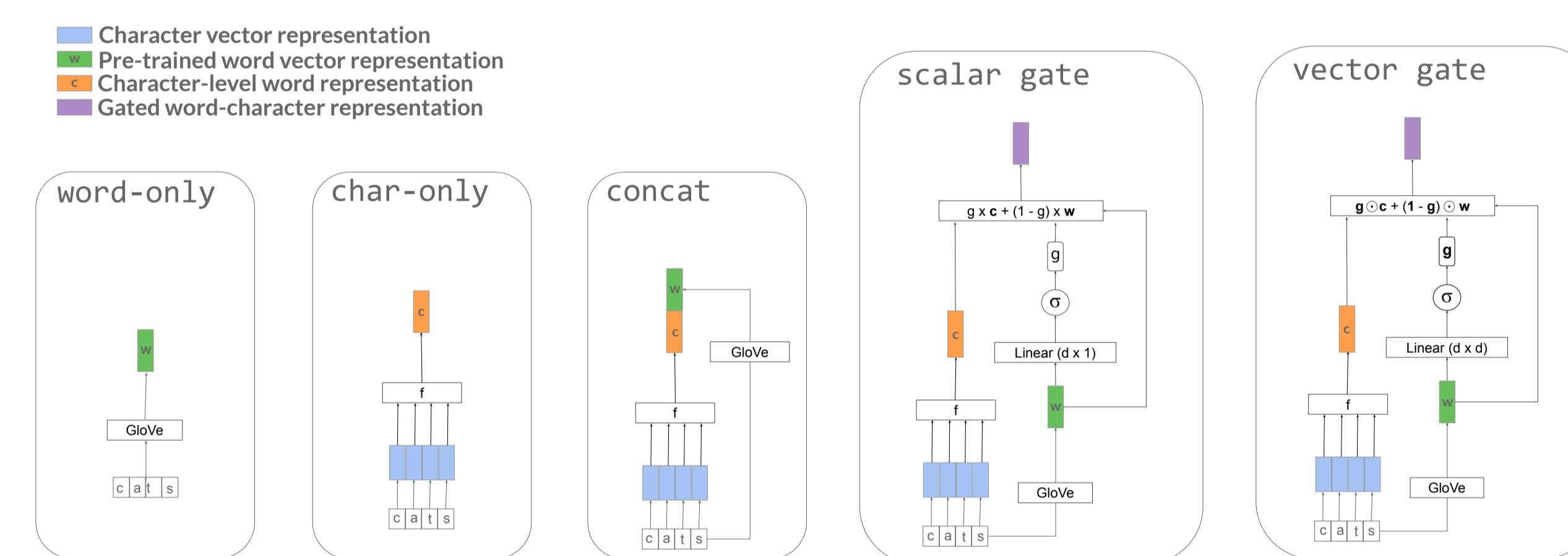
## Summary

- A vector gate is the best at combining character and word representations, as measured by word similarity tasks.
- This mechanism learns that to properly model increasingly infrequent words, it has to increasingly rely on character-level information.
- Despite the increased expressivity of word representations it offers, it has no clear effect in sentence representations, as measured by sentence evaluation tasks.

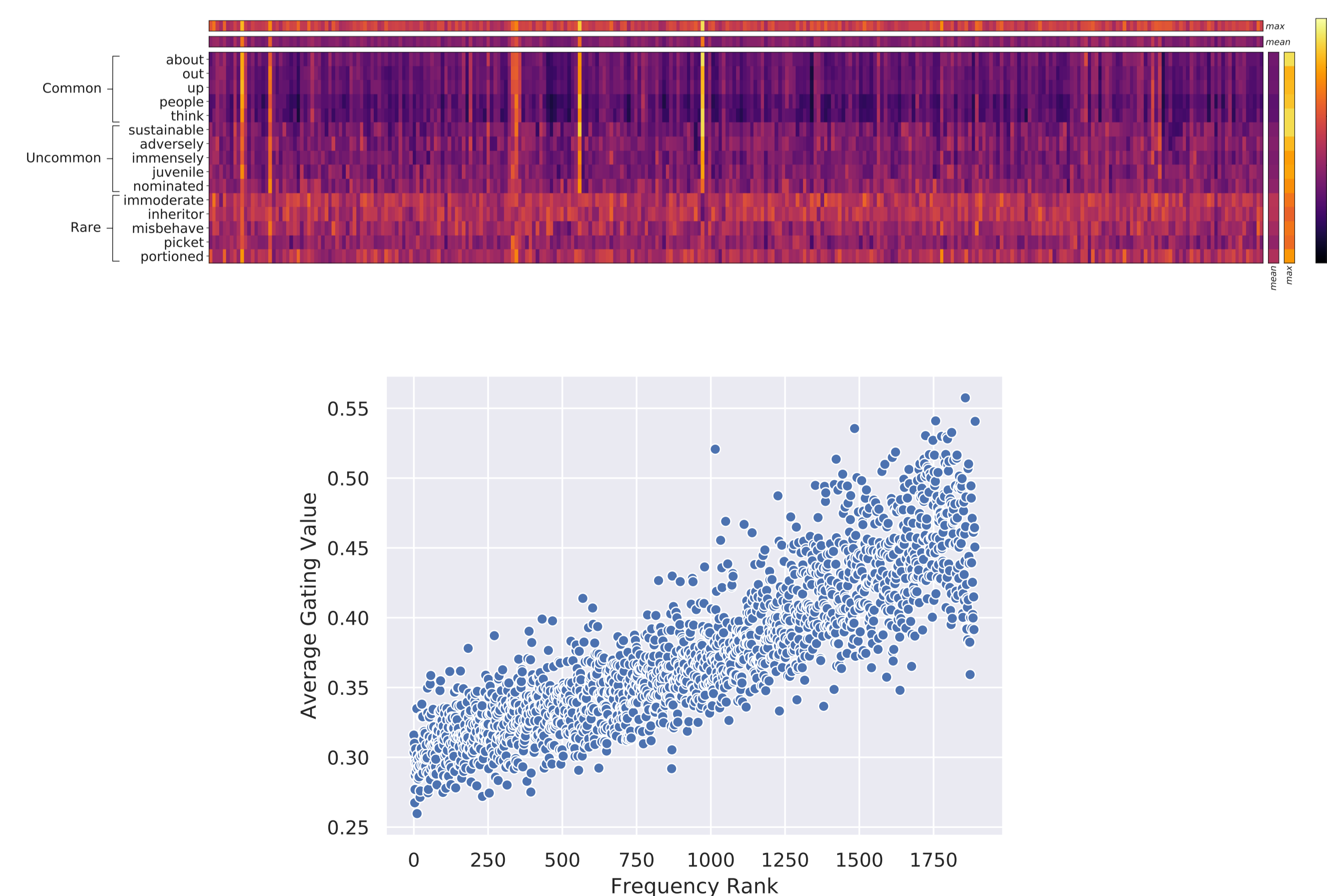
pool its output to obtain a sentence representation.

- Models are trained in the SNLI and MultiNLI (MNLI) datasets. The combined representation of the premise-hypothesis pair is defined as  $s = [s_p; s_h; |s_p - s_h|; s_p \odot s_h]$ , where  $s_p$  is the sentence representation of the premise and  $s_h$  that of the hypothesis.
- $s$  is finally mapped to label space through a fully-connected network.

## Gating Mechanisms



## Gating Values



## Model Description

- We initialize word-level word representations  $v_i^{(w)}$  with GloVe, and create character-level word representations  $v_i^{(c)}$  with a BiLSTM over randomly-initialized character representations.
- We test 5 different ways of combining  $v_i^{(w)}$  and  $v_i^{(c)}$  into the final word representations  $v_i$ :

• **scalar gate:**

$$g_i = \sigma(w^T v_i^{(w)} + b)$$

$$v_i = g_i v_i^{(c)} + (1 - g_i) v_i^{(w)}$$

• **vector gate:**

$$g_i = \sigma(W v_i^{(w)} + b)$$

$$v_i = g_i \odot v_i^{(c)} + (1 - g_i) \odot v_i^{(w)}$$

- **word-only:**  $v_i = v_i^{(w)}$
- **char-only:**  $v_i = v_i^{(c)}$
- **concat:**  $v_i = [v_i^{(w)}; v_i^{(c)}]$

- We feed the final word representations  $v_i$  to a BiLSTM, and max-

## Word-level Results

		MEN	MTurk287	MTurk771	RG65	RW	SimLex999	SimVerb3500	WS353	WS353R	WS353S
SNLI	w	71.78	35.40	49.05	61.80	18.43	19.17	10.32	39.27	28.01	53.42
	c	9.85	-5.65	0.82	-5.28	17.81	0.86	2.76	-2.20	0.20	-3.87
	cat	71.91	<b>35.52</b>	48.84	62.12	18.46	19.10	10.21	39.35	28.16	53.40
	sg	70.49	34.49	46.15	59.75	18.24	17.20	8.73	35.86	23.48	50.83
	vg	<b>80.00</b>	32.54	<b>62.09</b>	<b>68.90</b>	<b>20.76</b>	<b>37.70</b>	<b>20.45</b>	<b>54.72</b>	<b>47.24</b>	<b>65.60</b>
MNLI	w	68.76	50.15	68.81	65.83	18.43	42.21	25.18	61.10	58.21	70.17
	c	4.84	0.06	1.95	-0.06	12.18	3.01	1.52	-4.68	-3.63	-3.65
	cat	68.77	50.40	68.77	65.92	18.35	42.22	25.12	61.15	58.26	70.21
	sg	67.66	49.58	68.29	64.84	18.36	41.81	24.57	60.13	57.09	69.41
	vg	<b>76.69</b>	<b>56.06</b>	<b>70.13</b>	<b>69.00</b>	<b>25.35</b>	<b>48.40</b>	<b>35.12</b>	<b>68.91</b>	<b>64.70</b>	<b>77.23</b>

## Sentence-level Results

		Classification						Entailment	Relatedness	Semantic Textual Similarity		
		CR	MPQA	MR	SST2	SST5	SUBJ			TREC	SICKE	SICKR <sup>†</sup>
SNLI	w	80.50	84.59	74.18	78.86	42.33	<b>90.38</b>	<b>86.83</b>	86.37	88.52	59.90*	71.29*
	c	74.90*	78.86*	65.93*	69.42*	35.56*	82.97*	83.31*	84.13*	83.89*	59.33*	67.20*
	cat	80.44	84.66	74.31	78.37	41.34*	90.28	85.80*	<b>86.40</b>	88.44	59.90*	71.24*
	sg	<b>80.59</b>	84.60	<b>74.49</b>	<b>79.04</b>	41.63*	90.16	86.00	86.10*	<b>88.57</b>	60.05*	71.34*
	vg	80.42	<b>84.66</b>	74.26	78.87	<b>42.38</b>	90.07	85.97	85.67	<b>88.31*</b>	<b>60.92</b>	<b>71.99</b>
MNLI	w	83.80	<b>89.13</b>	79.05	83.38	45.21	91.79	89.23	84.92	86.33	66.08	71.96*
	c	70.23*	72.19*	62.83*	64.55*	32.47*	79.49*	74.74*	81.53*	75.92*	51.47*	61.74*
	cat	<b>83.96</b>	89.12	<b>79.23</b>	83.70	45.08*	<b>91.92</b>	<b>90.03</b>	<b>85.06</b>	86.45	<b>66.17</b>	71.82*
	sg	83.88	89.06	79.22	83.71	45.26	91.66*	88.83*	84.96	86.40	65.49*	71.87*
	vg	83.45*	89.05	79.13	<b>83.87</b>	<b>45.88</b>	91.55*	89.49	84.82	<b>86.50</b>	65.75	<b>72.82</b>

## Correlations Between Sentence and Word-level Tasks

