

Multi-Perspective Sentence Similarity Modeling with Convolutional Neural Networks

Hua He

University of Maryland College Park



Kevin Gimpel

Toyota Technological Institute at Chicago



Jimmy Lin

University of Waterloo



Problem: Sentence Similarity Measurement

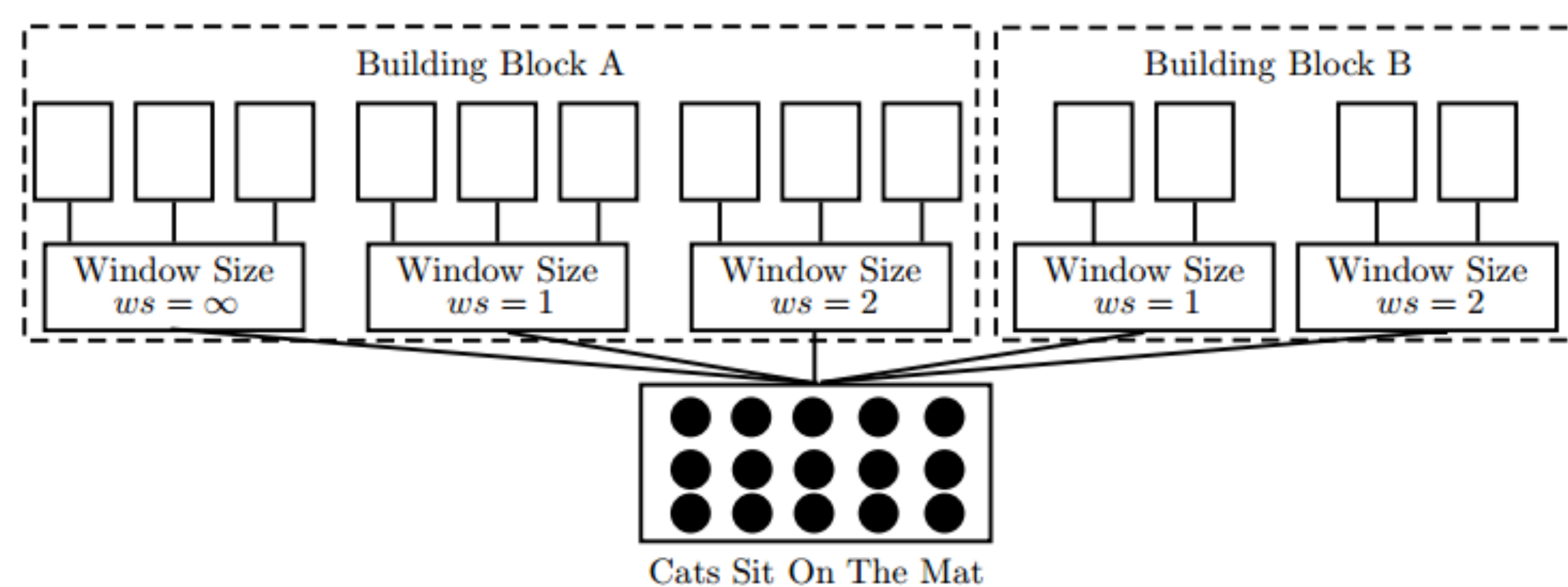
Given two sentences, measure their similarity:

The product also streams internet radio and comes with a 30-day free trial for realnetworks' rhapsody music subscription. The device plays internet radio streams and comes with a 30-day trial of realnetworks' rhapsody music service.

Approach: Multi-Perspective Sentence Representation and Structured Similarity Measurement

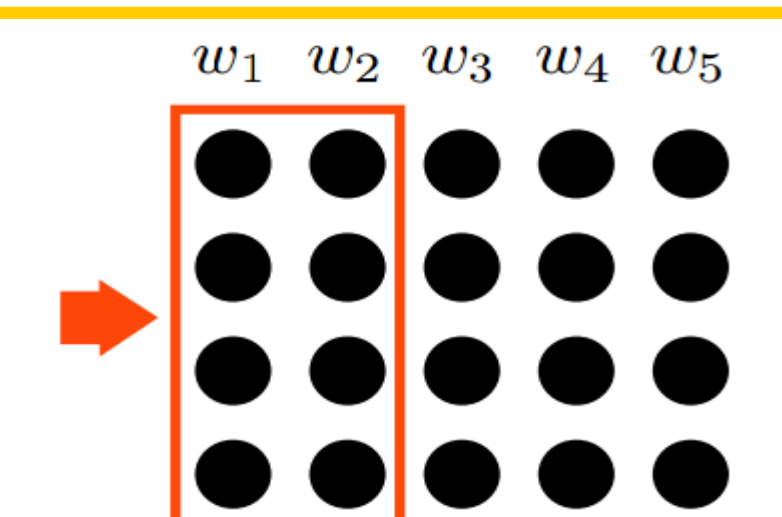
Part 1: Sentence Representation

to represent each sentence, we use **multiple types of convolution and pooling**

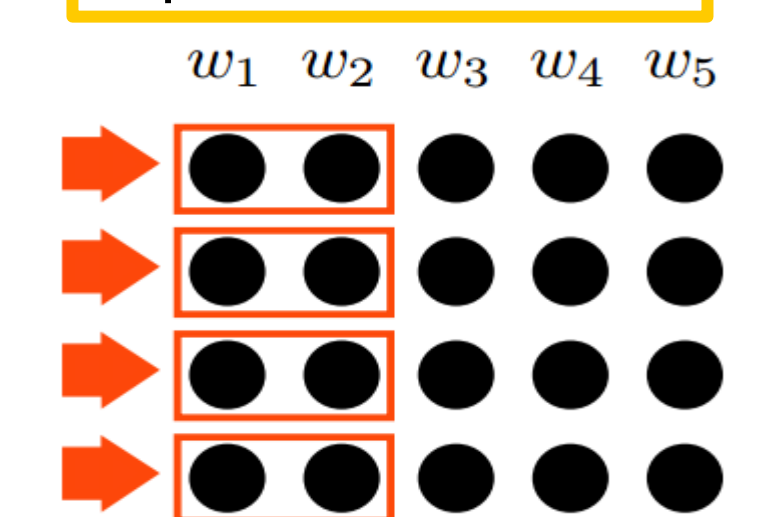


two types of convolution filters

Building Block A: holistic (all dimensions)

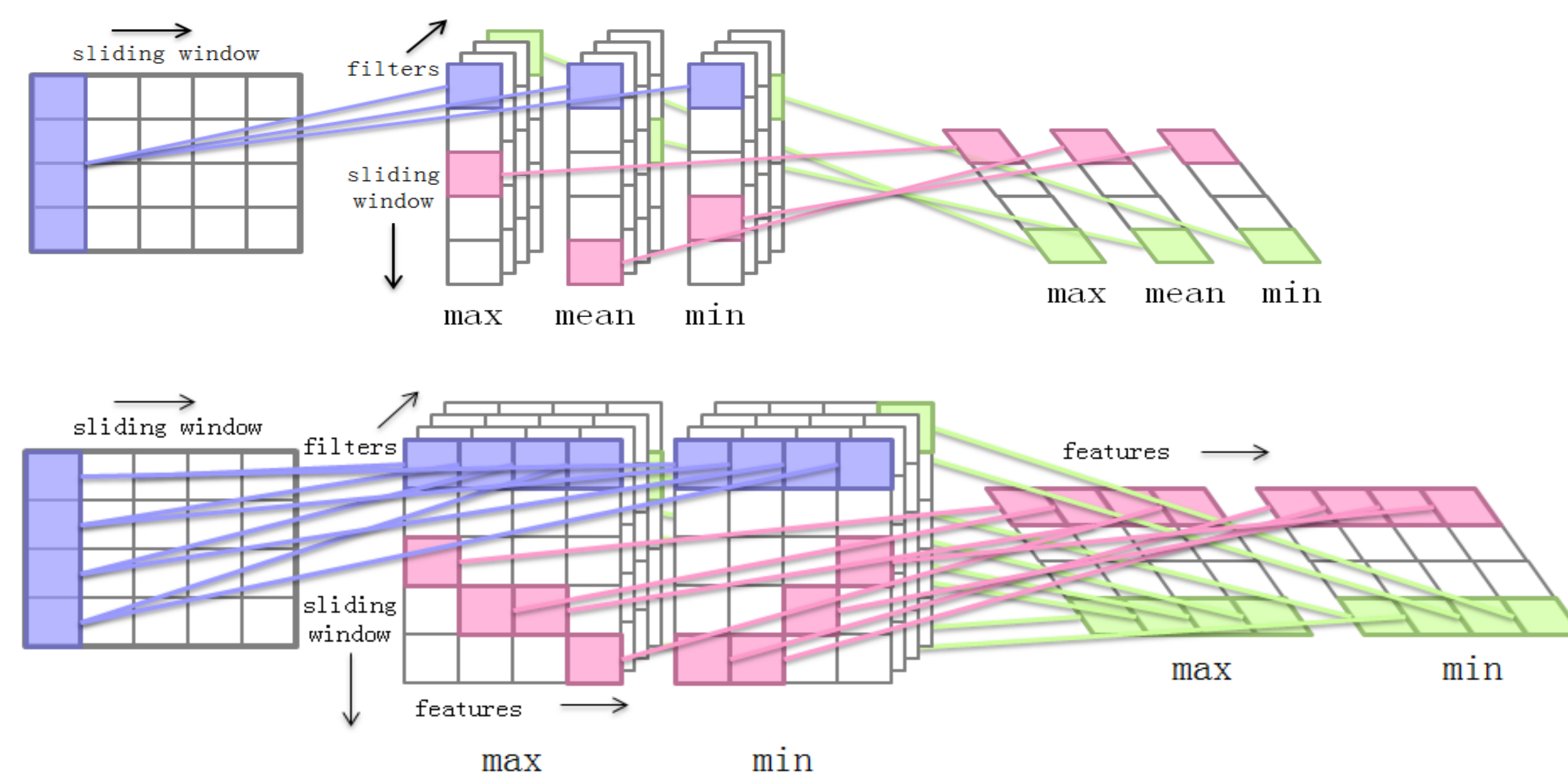


Building Block B: per-dimension



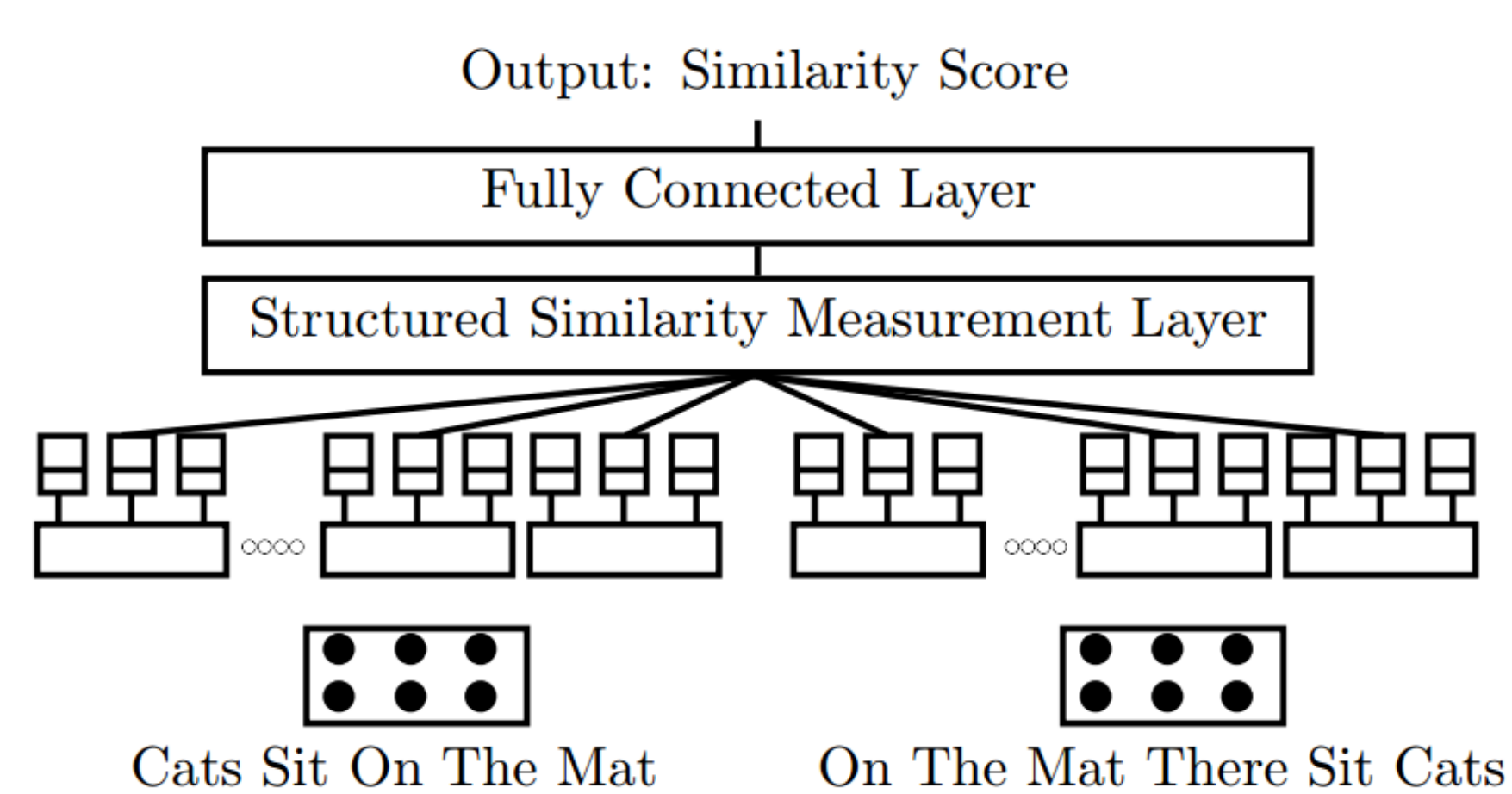
three types of pooling: max/min/mean

- each pooling group has **multiple window sizes** (1,2,3, infinity)
- each pooling group has independent underlying filters



Part 2: Structured Similarity Measurement

sentence representations compared by **structured similarity measurement layer**



two algorithms compare multiple pairs of local regions of sentence representations

Algorithm 1 Horizontal Comparison

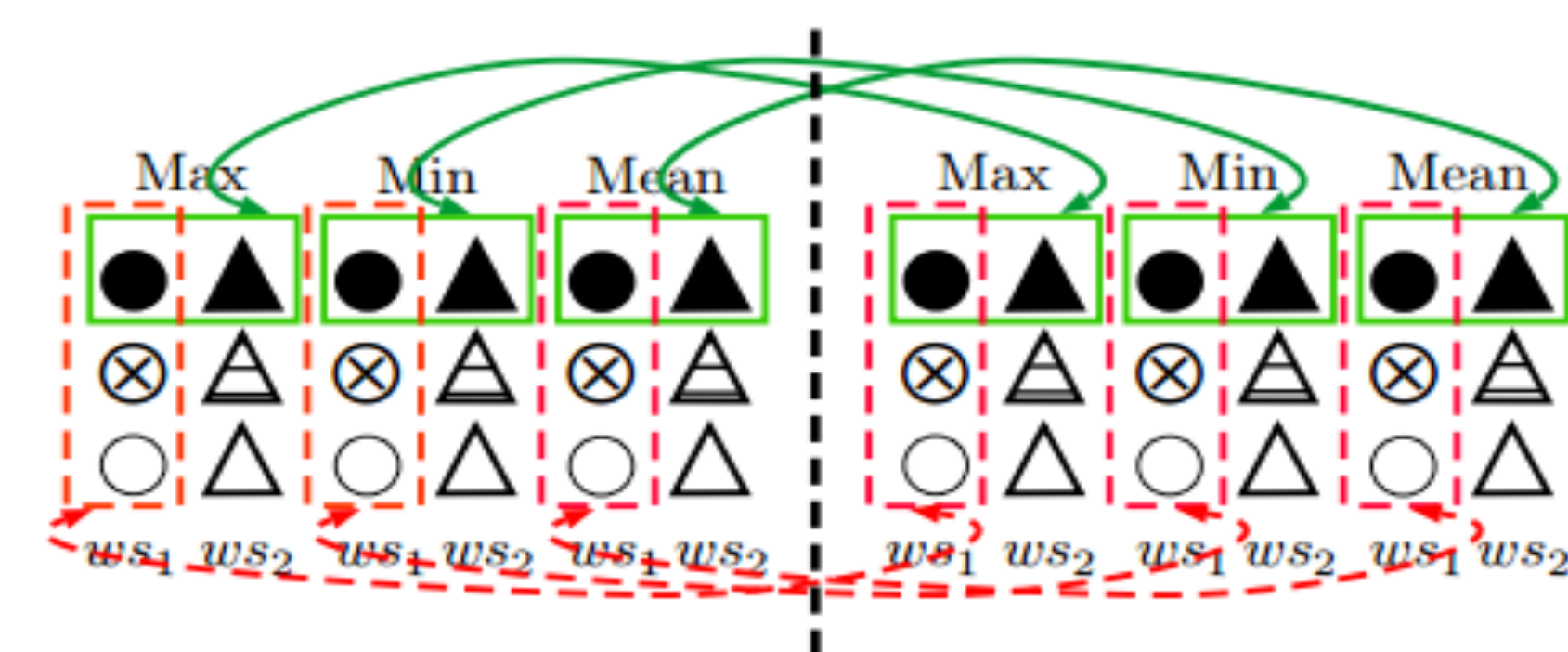
```

1: for each pooling p = max, min, mean do
2:   for each width ws1 = 1..n, infinity do
3:     regM1[*][ws1] = groupA(ws1, p, S1)
4:     regM2[*][ws1] = groupA(ws1, p, S2)
5:   end for
6:   for each i = 1..numFilterA do
7:     feaA = comU2(regM1[i], regM2[i])
8:     accumulate feaA for final layer
9:   end for
10: end for
    
```

Algorithm 2 Vertical Comparison

```

1: for each pooling p = max, min, mean do
2:   for each width ws1 = 1..n, infinity do
3:     oG1A = groupA(ws1, p, S1)
4:     for each width ws2 = 1..n, infinity do
5:       oG2A = groupA(ws2, p, S2)
6:       feaA = comU1(oG1A, oG2A)
7:       accumulate feaA for final layer
8:     end for
9:   end for
10:  for each width ws1 = 1..n do
11:    oG1B = groupB(ws1, p, S1)
12:    oG2B = groupB(ws1, p, S2)
13:    for each i = 1..numFilterB do
14:      feaB = comU1(oG1B[*][i], oG2B[*][i])
15:      accumulate feaB for final layer
16:    end for
17:  end for
18: end for
    
```



each uses **multiple similarity metrics** for vector comparison

$$comU_1(\vec{x}, \vec{y}) = \{\cos(\vec{x}, \vec{y}), L_2 Euclid(\vec{x}, \vec{y}), |\vec{x} - \vec{y}|\}$$

$$comU_2(\vec{x}, \vec{y}) = \{\cos(\vec{x}, \vec{y}), L_2 Euclid(\vec{x}, \vec{y})\}$$

simplified example of local region comparison for two sentences (Block A with 3 filters)

- green solid** lines: Horizontal Comparison (Alg. 1)
- red dotted** lines: Vertical Comparison (Alg. 2)

Experimental Results on Three Datasets

Experimental Setup

- Classification: Microsoft Research Paraphrase Corpus (MSRP)
- Similarity: Sentences Involving Compositional Knowledge (SICK)
- Similarity: Microsoft Video Paraphrase Corpus (MSRVID)
- multiple embeddings:
 - 300-dim GloVe (all tasks)
 - 200-dim POS (MSRP only)
 - 25-dim PARAGRAM (MSRP only)
- number of filters in Block A:
 - 525 (GloVe+POS+PARAGRAM) for MSRP
 - 300 for SICK/MSRVID
- embedding updating for MSRP only
- hinge loss for MSRP, KL-divergence loss (Tai et al., 2015) for SICK/MSRVID

MSRP

Model	Acc.	F1
Hu et al. (2014) ARC-I	69.6%	80.3%
Hu et al. (2014) ARC-II	69.9%	80.9%
Blacoe and Lapata (2012)	73.0%	82.3%
Fern and Stevenson (2008)	74.1%	82.4%
Finch (2005)	75.0%	82.7%
Das and Smith (2009)	76.1%	82.7%
Wan et al. (2006)	75.6%	83.0%
Socher et al. (2011)	76.8%	83.6%
Madnani et al. (2012)	77.4%	84.1%
Ji and Eisenstein (2013)	80.41%	85.96%
Yin and Schütze (2015) (without pretraining)	72.5%	81.4%
Yin and Schütze (2015) (with pretraining)	78.1%	84.4%
Yin and Schütze (2015) (pretraining+sparse features)	78.4%	84.6%
This work	78.60%	84.73%

SICK

Model	r	rho	MSE
Socher et al. (2014) DT-RNN	0.7863	0.7305	0.3983
Socher et al. (2014) SDT-RNN	0.7886	0.7280	0.3859
Lai and Hockenmaier (2014)	0.7993	0.7538	0.3692
Jimenez et al. (2014)	0.8070	0.7489	0.3550
Bjerva et al. (2014)	0.8268	0.7721	0.3224
Zhao et al. (2014)	0.8414	-	-
LSTM	0.8477	0.7921	0.2949
Bi-LSTM	0.8522	0.7952	0.2850
2-layer LSTM	0.8411	0.7849	0.2980
2-layer Bidirectional LSTM	0.8488	0.7926	0.2893
Tai et al. (2015) Const. LSTM	0.8491	0.7873	0.2852
Tai et al. (2015) Dep. LSTM	0.8676	0.8083	0.2532
This work	0.8686	0.8047	0.2606

MSRVID

Model	Pearson's r
Rios et al. (2012)	0.7060
Wang and Cer (2012)	0.8037
Beltagy et al. (2014)	0.8300
Bär et al. (2012)	0.8730
Šarić et al. (2012)	0.8803
This work	0.9090

Ablation Study

Ablation Component	MSRP Accuracy Diff.	MSRVID Pearson Diff.	SICK Pearson Diff.
Remove POS embeddings	-0.81	NA	NA
Remove PARAGRAM embeddings	-1.33	NA	NA
Remove per-dimension embeddings, building block A only	-0.75	-0.0067	-0.0014
Remove min and mean pooling, use max pooling only	-0.58	-0.0112	+0.0001
Remove multiple widths, ws = 1 and ws = infinity only	-2.14	-0.0048	-0.0012
Remove cosine and L2 Euclid distance in comU*	-2.31	-0.0188	-0.0309
Remove Horizontal Algorithm	-0.92	-0.0097	-0.0117
Remove Vertical Algorithm	-2.15	-0.0063	-0.0027
Remove similarity layer (completely flatten)	-1.90	-0.0121	-0.0288

Nine components in four groups:
 (1) input embeddings
 (2) sentence representation
 (3) similarity measurement metrics
 (4) similarity measurement layer

References

- J. Wieting, M. Bansal, K. Gimpel, K. Livescu, and D. Roth. 2015. From paraphrase database to compositional paraphrase model and back. *TACL*.
- K. S. Tai, R. Socher, and C. D. Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. *ACL*.
- W. Yin and H. Schütze. 2015. Convolutional neural network for paraphrase identification. *NAACL*.

code available: [hohocode.github.io/textSimilarityConvNet/](https://github.com/hohocode/textSimilarityConvNet/)