

# Capsule Networks for NLP

Will Merrill  
Advanced NLP  
10/25/18

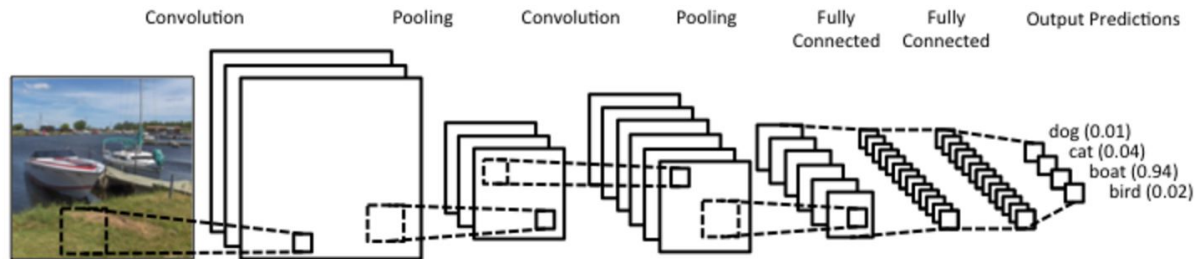
# Capsule Networks: A Better ConvNet

- Architecture proposed by Hinton as a replacement for ConvNets in computer vision
- Several recent papers applying them to NLP:
  - Zhao et al., 2018
  - Srivastava et al., 2018
  - Xia et al. 2018
- Goals:
  - Understand the architecture
  - Go through recent papers

What's Wrong with ConvNets?

# Convolutional Neural Networks

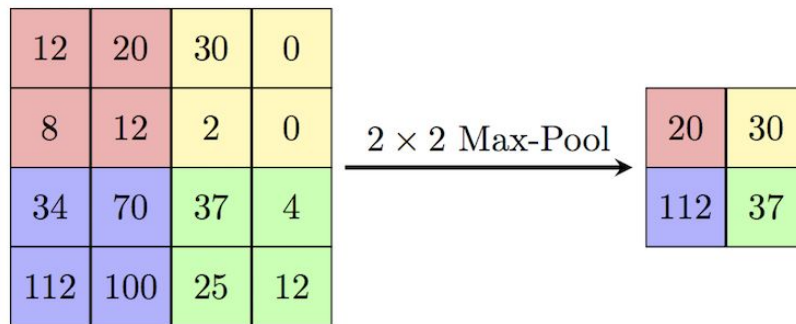
- Cascade of convolutional layers and max-pooling layers
- Convolutional layer:
  - Slide window over image and apply filter



<https://towardsdatascience.com/build-your-own-convolution-neural-network-in-5-mins-4217c2cf964f>

# Max-Pooling

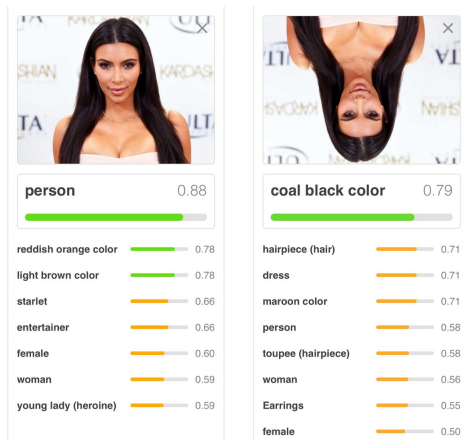
- ConvNets use max-pooling to move from low-level representations to high-level representations



[https://computersciencewiki.org/index.php/Max-pooling\\_-\\_Pooling](https://computersciencewiki.org/index.php/Max-pooling_-_Pooling)

# Problem #1: Transformational Invariance

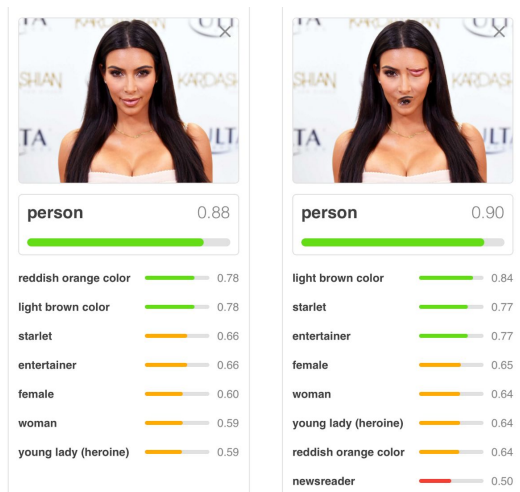
- We would like networks to recognize transformations of the same image
- Requires huge datasets of transformed images to learn transformations of high-level features



<https://medium.freecodecamp.org/understanding-capsule-networks-ais-alluring-new-architecture-bdb228173ddc>

# Problem #2: Feature Agreement

- Max-pooling in images loses information about relative position
- More abstractly, lower level features do not need to “agree”



<https://medium.freecodecamp.org/understanding-capsule-networks-ais-alluring-new-architecture-bdb228173ddc>

# Capsule Network Architecture

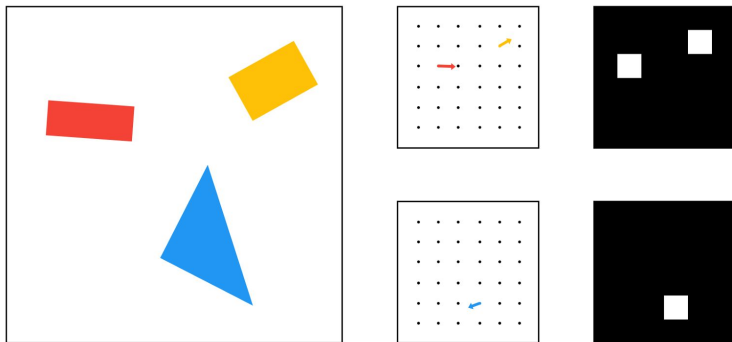


# Motivation

- We can solve problems #1 and #2 by attaching “instantiation parameters” to each filter
  - ConvNet: Is there a house here?
  - CapsNet: Is there a house with width  $w$  and rotation  $r$  here?
- Each filter at each position has a vector value instead of a scalar
- This vector is called a capsule

# Capsules

- The value of capsule  $i$  at some position is a vector  $\mathbf{u}_i$
- $|\mathbf{u}_i| \in (0, 1)$  gives the probability of existence of feature  $i$
- Direction of  $\mathbf{u}_i$  encodes the instantiation parameters of feature  $i$



# Capsules (Continued)



# Capsule Squashing Function

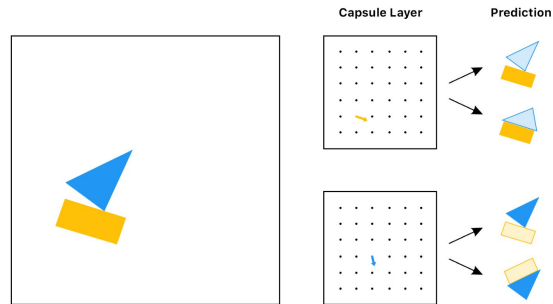
- New squashing function which puts magnitude of vector into (0, 1)
- Referred to in literature as  $g(..)$  or  $\text{squash}(..)$
- Will be useful later on

$$\mathbf{v}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|}$$

Sabour et al., 2017

# Routing by Agreement

- Capture child-parent relationships
- Combine features into higher-level ones only if the lower-level features “agree” locally
- Is this picture a house or a sailboat?



# Routing: Vote Vectors

- Learned transformation for what information should be “passed up” to the next layer
- Models what information is relevant for abstraction/agreement
- $\hat{u}_{j|i}$  denotes the vote vector from capsule  $i$  to capsule  $j$  in the next layer

$$\hat{u}_{j|i} = W_j^{c1} u_i + \hat{b}_{j|i}$$

Zhao et al., 2018

# Routing: Dynamic Routing Algorithm

- Unsupervised iterative method for computing routing
- No parameters (But depends on vote vectors)
- Used to connect capsule layers
- Compute next layer of capsules  $\{\mathbf{v}_j\}$  from vote vectors

---

**Procedure 1** Routing algorithm.

---

```
1: procedure ROUTING( $\hat{\mathbf{u}}_{j|i}, r, l$ )
2:   for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow 0$ .
3:   for  $r$  iterations do
4:     for all capsule  $i$  in layer  $l$ :  $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$  ▷ softmax computes Eq. 3
5:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}$ 
6:     for all capsule  $j$  in layer  $(l + 1)$ :  $\mathbf{v}_j \leftarrow \text{squash}(\mathbf{s}_j)$  ▷ squash computes Eq. 1
7:     for all capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $(l + 1)$ :  $b_{ij} \leftarrow b_{ij} + \hat{\mathbf{u}}_{j|i} \cdot \mathbf{v}_j$ 
   return  $\mathbf{v}_j$ 
```

---

# Types of Capsule Layers

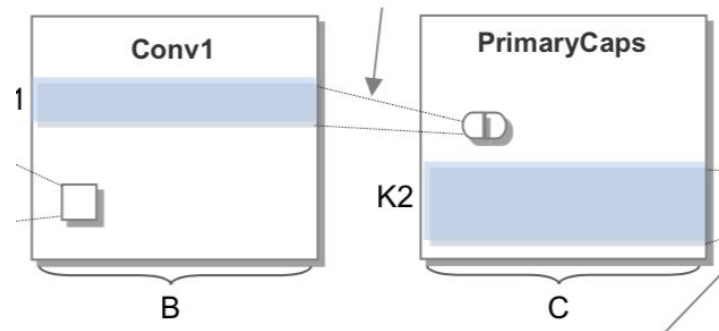
1. **Primary Capsule Layer:** Convolutional output → capsules
2. **Convolutional Capsule Layer:** Local capsules → capsules
3. **Feedforward Capsule Layer:** All capsules → capsules



# Primary Capsule Layer

Convolutional output  $\rightarrow$  capsules

Create  $C$  capsules from  $B$  filters

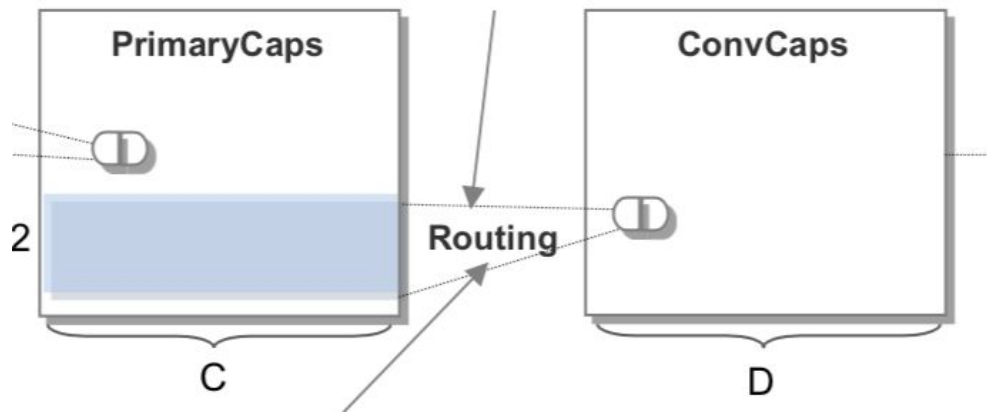


1. Convolution output with  $B$  filters:  $\mathbf{M} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_B] \in \mathbb{R}^{(L-K_1+1) \times B}$
2. Transform each row of features:  $p_i = g(W^b \mathbf{M}_i + \mathbf{b}_1) \quad W^b \in \mathbb{R}^{B \times d}$
3. Collect  $C$   $d$ -dimensional capsules:  $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_C] \in \mathbb{R}^{(L-K_1+1) \times C \times d}$

# Convolutional Capsule Layer

Local capsules in layer #1  $\rightarrow$  capsules in layer #2

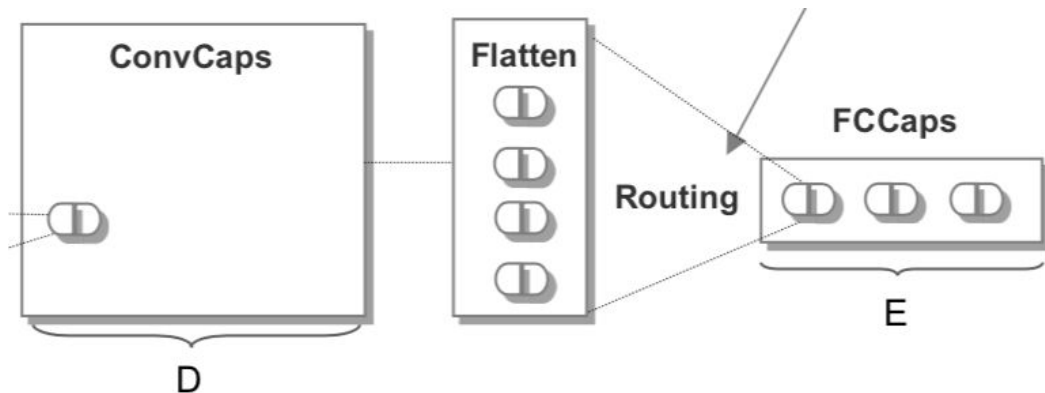
- Route a sliding window of capsules in previous layer into capsules in next layer



# Feedforward Capsules Layer

All capsules in layer #1  $\rightarrow$  capsules in layer #2

1. Flatten all capsules in layer #1 into a vector
2. Route from this vector of capsules into new capsules



# Margin Loss

- Identify each output capsule with a class
- Classification loss for capsules
- Calculate on output of feedforward capsule layer
- Ensures that the capsule vector for the correct class is long ( $||\mathbf{v}|| \approx 1$ )

$$L_k = T_k \max(0, m^+ - ||\mathbf{v}_k||)^2 + \lambda (1 - T_k) \max(0, ||\mathbf{v}_k|| - m^-)^2$$

Sabour et al., 2017

# Investigating Capsule Networks with Dynamic Routing for Text Classification

Zhao, Ye, Yang, Lei, Zhang, Zhao 2018

# Main Ideas

1. Develops capsule network architecture for text classification tasks
2. Achieves state-of-the-art performance on single-class text classification
3. Capsules allow transferring single-class classification knowledge to multi-class task very well

# Text Classification

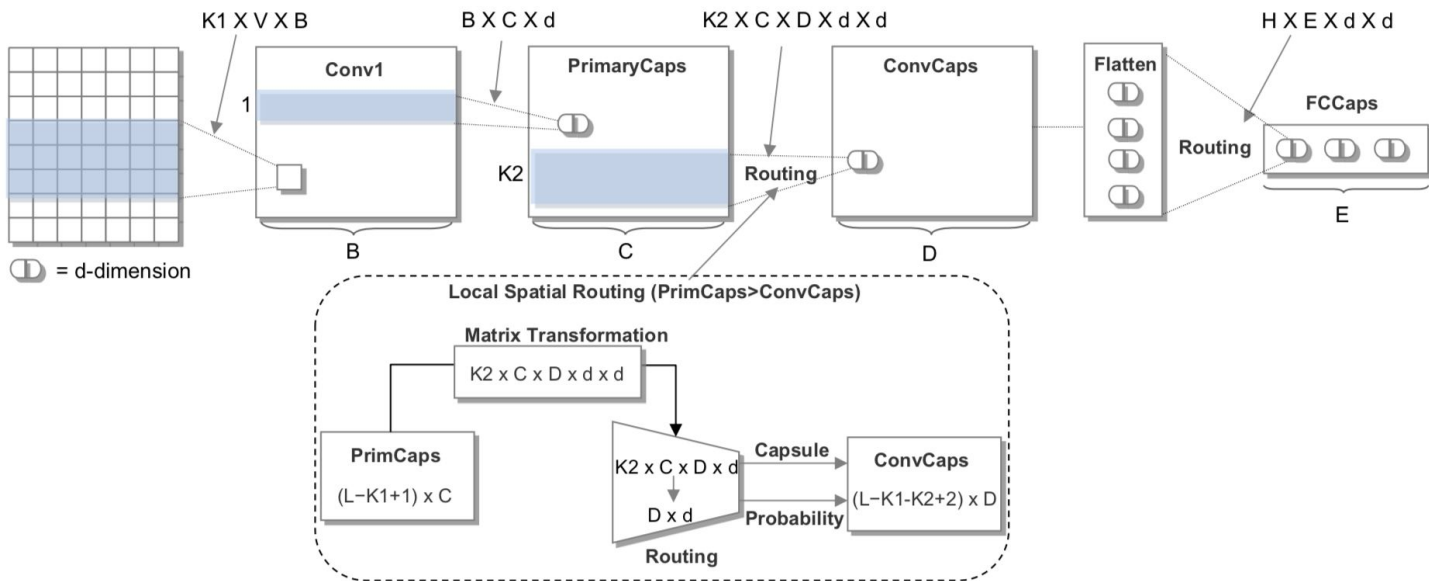
- Read text and classify something about the passage
- Sentiment analysis, toxicity detection, etc.

# Multi-Class Text Classification

- Document can be labeled as multiple classes
  - Example: In toxicity detection, Toxic and Threatening

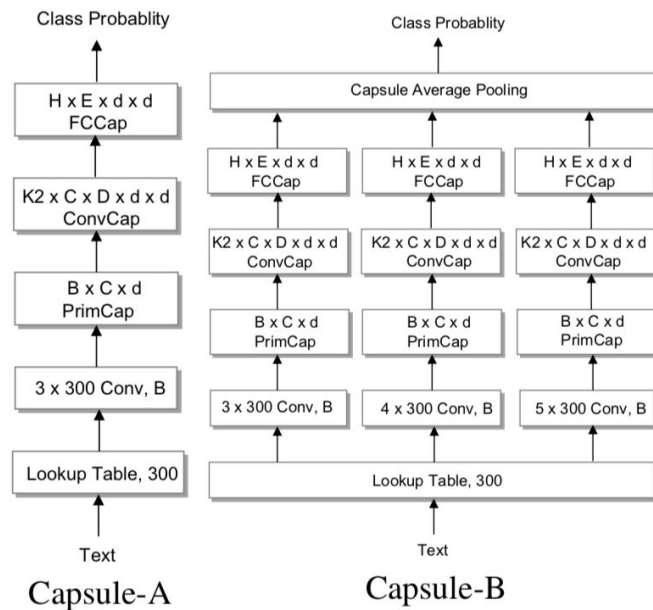


# Text Classification Architecture



# Architectural Variants

- **Capsule-A:** One capsule network
- **Capsule-B:** Three capsule networks that are averaged at the end



# Orphan Category

- Add a capsule that corresponds to no class to the final layer
- Network can send words unimportant to classification to this category
  - Function words like *the*, *a*, *in*, etc.
- More relevant in the NLP domain than in images because images don't have a “default background”

# Datasets

Single-Label

Dataset	Train	Dev	Test	Classes	Classification Task
MR	8.6k	0.9k	1.1k	2	review classification
SST-2	8.6k	0.9k	1.8k	2	sentiment analysis
Subj	8.1k	0.9k	1.0k	2	opinion classification
TREC	5.4k	0.5k	0.5k	6	question categorization
CR	3.1k	0.3k	0.4k	2	review classification
AG's news	108k	12.0k	7.6k	4	news categorization

Multi-Label

Dataset	Train	Dev	Test	Description
Reuters-Multi-label	5.8k	0.6k	0.3k	only multi-label data in test
Reuters-Full	5.8k	0.6k	3.4k	full data in test

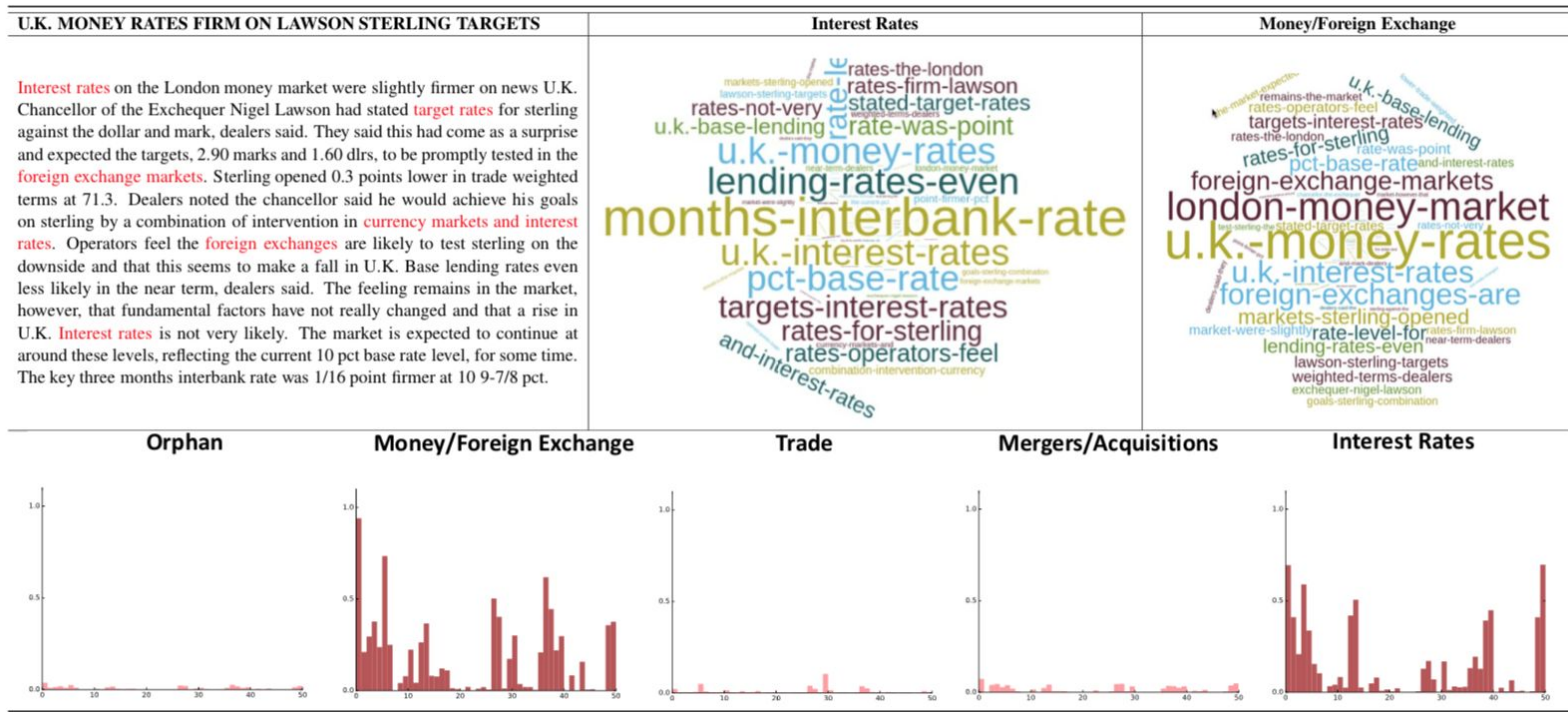
# Single-Class Results

	MR	SST2	Subj	TREC	CR	AG's
LSTM	75.9	80.6	89.3	86.8	78.4	86.1
BiLSTM	79.3	83.2	90.5	89.6	82.1	88.2
Tree-LSTM	80.7	85.7	91.3	91.8	83.2	90.1
LR-LSTM	81.5	<b>87.5</b>	89.9	-	82.5	-
CNN-rand	76.1	82.7	89.6	91.2	79.8	92.2
CNN-static	81.0	86.8	93.0	92.8	84.7	91.4
CNN-non-static	81.5	87.2	93.4	<b>93.6</b>	84.3	92.3
CL-CNN	-	-	88.4	85.7	-	92.3
VD-CNN	-	-	88.2	85.4	-	91.3
Capsule-A	81.3	86.4	93.3	91.8	83.8	92.1
Capsule-B	<b>82.3</b>	86.8	<b>93.8</b>	92.8	<b>85.1</b>	<b>92.6</b>

# Multi-Class Transfer Learning Results

	Reuters-Multi-label				Reuters-Full			
	ER	Precision	Recall	F1	ER	Precision	Recall	F1
LSTM	23.3	86.7	54.7	63.5	62.5	78.6	72.6	74.0
BiLSTM	26.4	82.3	55.9	64.6	65.8	83.7	75.4	77.8
CNN-rand	22.5	88.6	56.4	67.1	63.4	78.7	71.5	73.6
CNN-static	27.1	91.1	59.1	69.7	63.3	78.5	71.2	73.3
CNN-non-static	27.4	92.0	59.7	70.4	64.1	80.6	72.7	75.0
Capsule-A	57.2	88.2	80.1	82.0	66.0	83.9	<b>80.5</b>	80.2
Capsule-B	<b>60.3</b>	<b>95.4</b>	<b>82.0</b>	<b>85.8</b>	<b>67.7</b>	<b>86.4</b>	80.1	<b>81.4</b>

# Connection Strength Visualization



# Discussion

- Capsule network performs strongly on single-class text-classification
- Capsule model transfers effectively from single-class to multi-class domain
  - Richer representation
  - No softmax in last layer
- Useful because multi-class data sets are hard to construct (exponentially larger than single-class data sets)



# Identifying Aggression and Toxicity in Comments Using Capsule Networks

Srivastava, Khurana, Tewari 2018

# Main Ideas

1. Develop end-to-end capsule model that outperforms state-of-the-art models for toxicity detection
2. Eliminate need for pipelining and preprocessing
3. Performs especially well on code-mixed comments (comments switching between English and Hindi)

# Toxicity Detection

- Human moderation of online content is expensive – useful to do algorithmically
- Classify comments as **toxic**, **severe toxic**, **identity hate**, etc.

◆ 92% similar to comments people said were  
"toxic"

SEEM WRONG?

Nieman Lab is a great website — only an idiot like you would think some other website could possibly be better. You dumb jerk.

● 2% similar to comments  
people said were "toxic"

SEEM WRONG?

I respectfully disagree

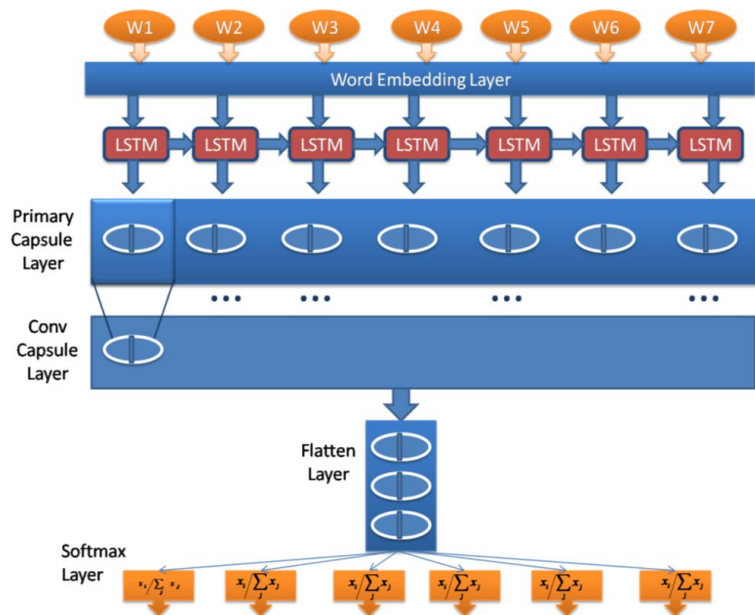
# Challenges in Toxicity Detection

- Out-of-vocabulary words
- Code-mixing of languages
- Class imbalance

# Why Capsule Networks?

- Seem to be good at text classification (Zhao et al., 2018)
- Should be better at code-mixing than sequential models (build up local representations)

# Architecture



- Very similar to architecture to Zhao et al.
- Feature extraction convolutional layer replaced by LSTM
- Standard softmax layer instead of margin loss

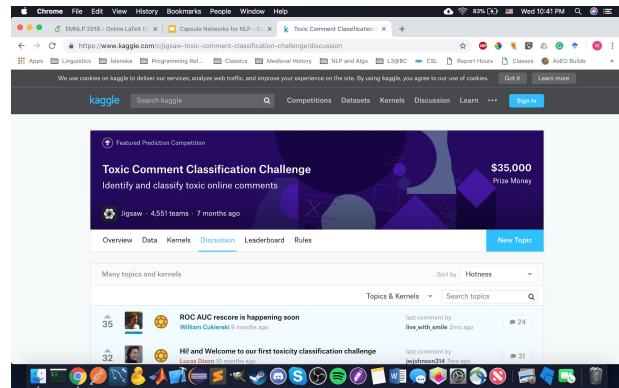
# Focal Loss

- Loss function on standard softmax output
- Used to solve the class imbalance problem
- Weights rare classes higher than cross-entropy

$$FL(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t), \text{ where } p_t = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{else} \end{cases}$$

# Datasets

- **Kaggle Toxic Comment Classification**
  - English
  - Classes: Toxic, Severe Toxic, Obscene, Threat, Insult, Identity Hate
- **First Shared Task on Aggression Identification (TRAC)**
  - Mixed English and Hindi
  - Classes: Overtly Aggressive, Covertly Aggressive, Non-Aggressive



<https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/discussion>

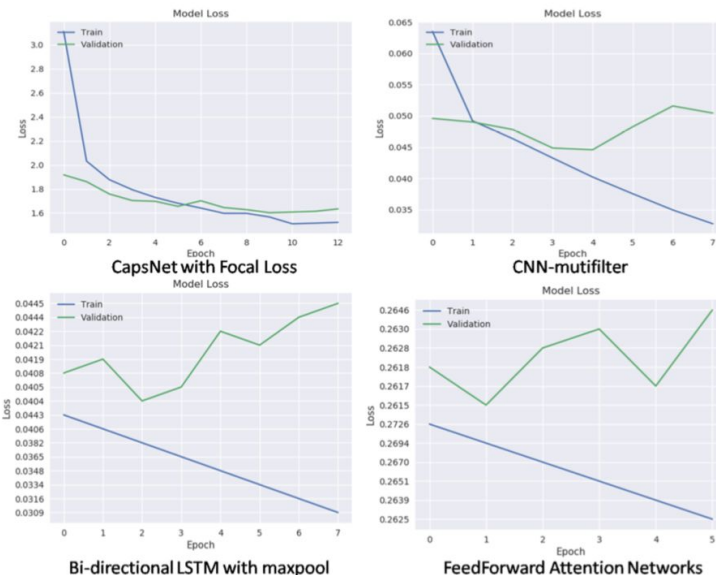


# Results

Model_Name	Kaggle-toxic comment classification (ROC-AUC)	TRAC - 1 (English-FB) (Weighted F1)	TRAC - 1 (English-TW) (Weighted F1)
<b>CNN-multifilter</b>	95.16	55.43	53.41
<b>CNN-LSTM</b>	96.85	62.20	47.68
<b>Bi-directional LSTM with maxpool</b>	97.35	59.79	51.146
<b>FeedForward Attention Networks</b>	97.42	57.43	55.49
<b>Hierarchical ConvNets</b>	97.95	51.38	50.43
<b>Bi-LSTM, Logistic Regression</b>	98.17	57.17	52.1
<b>Bi-LSTM, xgboosted</b>	98.19	57.33	52.31
<b>Bi-LSTM with skip connections</b>	98.20	61.78	51.98
<b>Pre-trained LSTMs</b>	98.25	60.18	58.7
<b>CapsuleNet without Focal Loss</b>	98.21	62.032	58.600
<b>CapsuleNet with Focal Loss</b>	<b>98.46</b>	<b>63.43</b>	<b>59.41</b>

# Training/Validation Loss

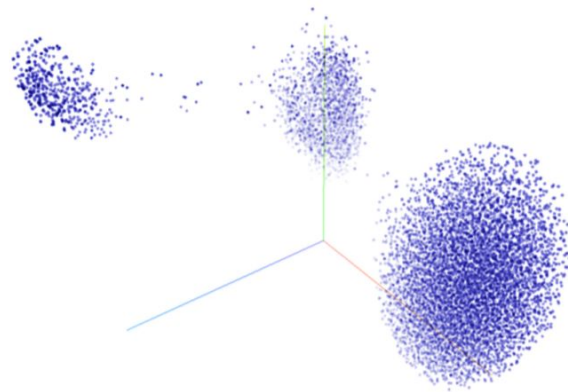
- Training and validation loss stayed much closer for the capsule model
- $\Rightarrow$  Avoids overfitting



(a) Training and Validation Loss for Kaggle Toxic Comment Classification Dataset

# Word Embeddings on Kaggle Corpus

- Three clear clusters:
  - Neutral words
  - Abusive words
  - Toxic words + place names



(b) Clusters for word obtained after training

# OOV Embeddings

NN to “politics”	NN to “bharat”
politic	bharatiya
politican	bhar
politico	mahabharata
politicize	bharti
politician	bhaskar

NN to “kut*e”(Hindi)
chu**ya
sa*le
tere
g**d
ma***rc**d

Table 2: Example of handling misspelt words and transliteration. NN : Nearest Neighbour

- Out of vocabulary words randomly initialized
- Converge to accurate vectors

# Discussion

- The novel capsule network architecture performed the best on all three datasets
- No data preprocessing done
- Avoids overfitting
- Local representations lead to big gains in mixed-language case

# Zero-shot User Intent Detection via Capsule Neural Networks

Xia, Zhang, Yan, Chang, Yu 2018

# Main Ideas

1. Capsule networks extract and organize information during supervised intent detection
2. These learned representations can be effectively transferred to the task of zero-shot intent detection

# User Intent Detection

- Text classification task for question answering and dialog systems
- Classify which action a user query represents out of a known set of actions
  - **GetWeather, PlayMusic**



# Zero-Shot User Intent Detection

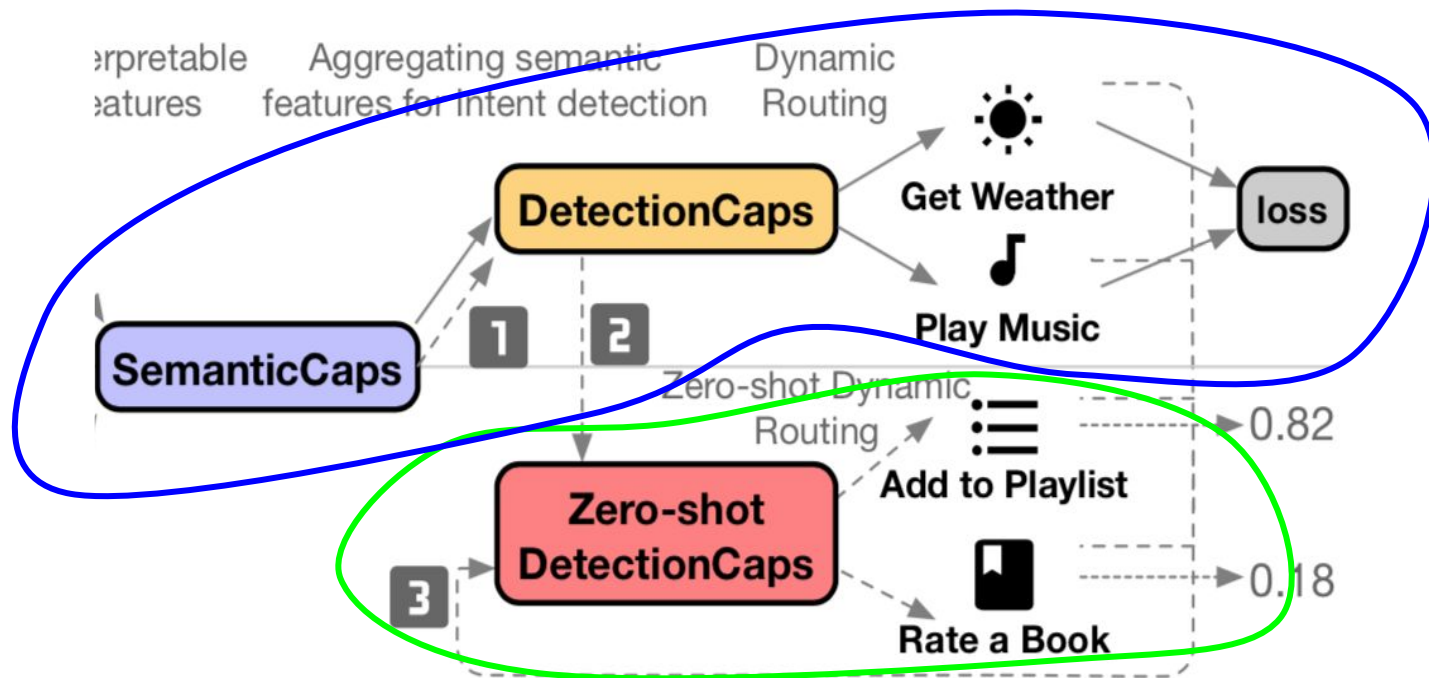
- Training set with known set of intents
  - **GetWeather, PlayMusic**
- Test set has unseen “emerging” intents
  - **AddToPlaylist, RateABook**
- Transfer information about known intents to new domain of emerging intents

# What Signal is There?

- Embedding of the string name of the unknown and known intents
- Output capsules for known intents
- Can combine these two things to do zero-shot learning

# Architecture

Network trained on known intents



Extension for zero-shot inference

# SemanticCaps Layer

- Extract features using self-attention LSTM

Combine to get  $\mathbf{H}$

$$\left\{ \begin{array}{l} \vec{\mathbf{h}}_t = \text{LSTM}_{fw}(\mathbf{w}_t, \vec{\mathbf{h}}_{t-1}), \\ \overleftarrow{\mathbf{h}}_t = \text{LSTM}_{bw}(\mathbf{w}_t, \overleftarrow{\mathbf{h}}_{t+1}). \end{array} \right\}$$

Self-attention weights

$$\left\{ \mathbf{A} = \text{softmax} \left( \mathbf{W}_{s2} \tanh \left( \mathbf{W}_{s1} \mathbf{H}^T \right) \right) \right\}$$

$\mathbf{M}$  is the extracted features

$$\left\{ \begin{array}{l} \mathbf{M} = \mathbf{A}\mathbf{H} \\ (\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_R) \in \mathbb{R}^{R \times 2D_H} \end{array} \right\}$$

# DetectionCaps Layer

- Standard convolutional capsule layer → feedforward capsule layer

# Loss During Training

- Normal max-margin loss + regularization
- Regularization incentivizes semantic capsules to capture different features
- Regularization controlled by hyperparameter  $\alpha$

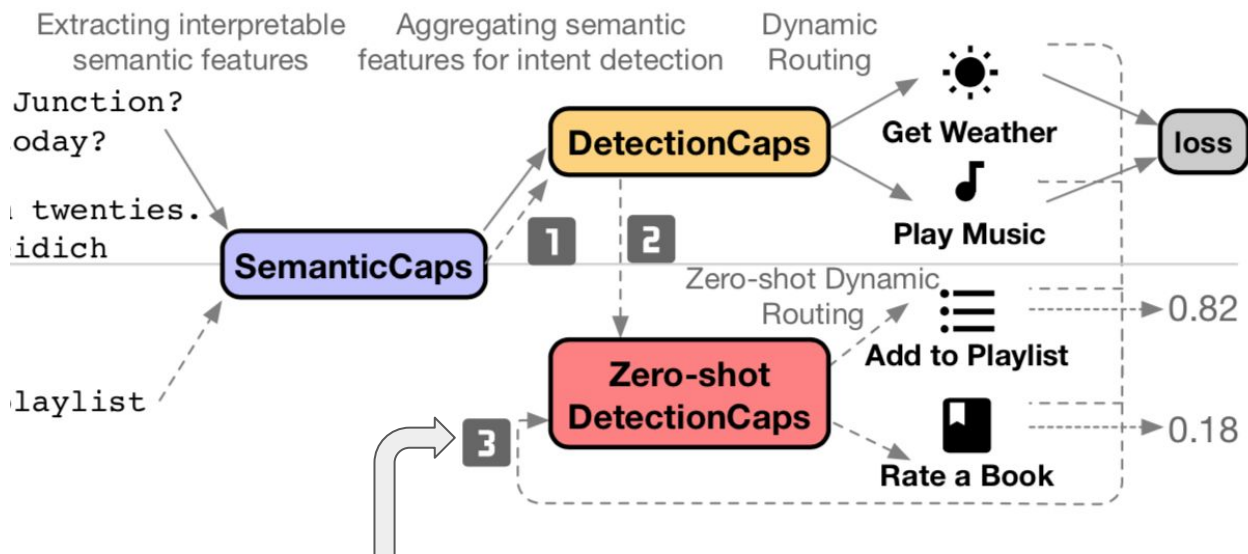
$$\begin{array}{l} \text{Max-margin loss} \\ \text{Regularization term} \end{array} \left\{ \begin{array}{l} \mathcal{L} = \sum_{k=1}^K \{ \mathbb{I}[y = y_k] \cdot \max(0, m^+ - \|\mathbf{v}_k\|)^2 \\ \quad + \lambda \mathbb{I}[y \neq y_k] \cdot \max(0, \|\mathbf{v}_k\| - m^-)^2 \} \\ \quad + \alpha \|\mathbf{A}\mathbf{A}^T - I\|_F^2, \end{array} \right\}$$

# Intent Detection Results

Model	SNIPS-NLU (on 5 existing intents)				CVA (on 80 existing intents)			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
TFIDF-LR	0.9546	0.9551	0.9546	0.9545	0.7979	0.8104	0.7979	0.7933
TFIDF-SVM	0.9584	0.9586	0.9584	0.9581	0.7989	0.8111	0.7989	0.7942
CNN	0.9595	0.9596	0.9595	0.9595	0.8223	0.8288	0.8223	0.8210
RNN	0.9516	0.9522	0.9516	0.9518	0.8286	0.8330	0.8286	0.8275
GRU	0.9535	0.9535	0.9535	0.9534	0.8239	0.8281	0.8239	0.8216
LSTM	0.9569	0.9573	0.9569	0.9569	0.8319	0.8387	0.8319	0.8306
Bi-LSTM	0.9501	0.9502	0.9501	0.9502	0.8428	0.8479	0.8428	0.8419
Self-attention Bi-LSTM	0.9524	0.9522	0.9524	0.9522	0.8521	0.8590	0.8521	0.8513
INTENTCAPSNET	<b>0.9621</b>	<b>0.9620</b>	<b>0.9621</b>	<b>0.9620</b>	<b>0.9088</b>	<b>0.9160</b>	<b>0.9088</b>	<b>0.9023</b>

# Architecture Revisited

- Goal: Use predicted capsules for known intents for zero-shot inference





# Generalizing to Emerging Intents

- Build similarity matrix between existing intents and emerging intents based on embeddings for intent names:

$$q_{lk} = \frac{\exp \{-d(\mathbf{e}_{z_l}, \mathbf{e}_{y_k})\}}{\sum_{k=1}^K \exp \{-d(\mathbf{e}_{z_l}, \mathbf{e}_{y_k})\}},$$

where

$$d(\mathbf{e}_{z_l}, \mathbf{e}_{y_k}) = (\mathbf{e}_{z_l} - \mathbf{e}_{y_k})^T \Sigma^{-1} (\mathbf{e}_{z_l} - \mathbf{e}_{y_k})$$

# Classifying Emerging Intents

1. Goal is to get prediction vector for emerging intent /
2. Have vote vectors  $\mathbf{g}_{k,r}$  from known intent classification
3. Represent vote vector for emerging intent as weighted sum of known intents:

$$\mathbf{u}_{l|r} = \sum_{k=1}^K q_{lk} \mathbf{g}_{k,r}$$

4. Use dynamic routing to get an activation capsule  $\mathbf{n}_l$  for each emerging intent
5. Pick the  $\mathbf{n}_l$  with largest magnitude

# Zero-Shot Intent Detection Results

Model	SNIPS-NLU (on 2 emerging intents)				CVA (on 20 emerging intents)			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
DeViSE (Frome et al., 2013)	0.7447	0.7448	0.7447	0.7446	0.7809	0.8060	0.7809	0.7617
CMT (Socher et al., 2013)	0.7396	<b>0.8266</b>	0.7396	0.7206	0.7721	0.7728	0.7721	0.7445
CDSSM (Chen et al., 2016a)	0.7588	0.7625	0.7588	0.7580	0.2140	0.4072	0.2140	0.1667
Zero-shot DNN (Kumar et al., 2017)	0.7165	0.7330	0.7165	0.7116	0.7903	0.8240	0.7903	0.7774
INTENTCAPSNET-ZSL w/o Self-attention	0.7587	0.7764	0.7588	0.7547	0.8103	0.8512	0.8103	0.8115
INTENTCAPSNET-ZSL w/o Bi-LSTM	0.7619	0.7631	0.7619	0.7616	0.8366	<b>0.8770</b>	0.8366	0.8403
INTENTCAPSNET-ZSL w/o Regularizer	0.7675	0.7676	0.7675	0.7675	0.8544	0.8730	0.8544	0.8553
INTENTCAPSNET-ZSL	<b>0.7752</b>	0.7762	<b>0.7752</b>	<b>0.7750</b>	<b>0.8628</b>	0.8751	<b>0.8629</b>	<b>0.8635</b>

# Discussion

- Representational power of capsule network can be leveraged for zero-shot learning
- Interesting regularizations and architectural extensions for capsule networks

# Conclusion

- Capsule representations encode “instantiation parameters” of features
- Papers follow a standard CapsNet architecture for text classification:
  - a. Features Extraction (ConvNet or LSTM)
  - b. Primary Capsule Layer
  - c. Convolutional Capsule Layer
  - d. Classification (Margin or softmax)
- Capsule representations can be leveraged for transfer/zero-shot learning

# Discussion Questions

1. What is powerful about capsule representations?
2. Are capsule networks good for NLP, or are they just good for vision?
3. Why has NLP capsule research focused on text classification tasks?
4. What are some other NLP tasks that capsule networks could be applied to?
5. What other advanced architectures could be useful in NLP?

# Other Papers

- Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. 2017. [Dynamic routing between capsules](#). In *Advances in Neural Information Processing Systems*, pages 3859–3869.

# Other Materials

- <https://medium.freecodecamp.org/understanding-capsule-networks-ais-alluring-new-architecture-bdb228173ddc>
- <https://medium.com/ai%C2%B3-theory-practice-business/understanding-hints-capsule-networks-part-i-intuition-b4b559d1159b>