

# Concept to Code:

## Aspect sentiment classification with Deep Learning

Muthusamy Chelliah  
Flipkart

Asif Ekbal  
IIT Patna

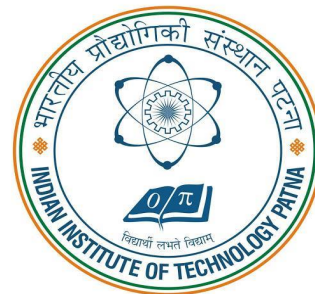
Mohit Gupta  
Flipkart

Tutorial at IJCAI-2019  
August 11, 2019, Macao



### Acknowledgement

Shad Akhtar, PhD Student, IIT Patna



# Outline

- Introduction (10 min.s) (Asif)
- LSTM/attention based ASC (20 min.s) (Asif)
- Code - LSTM/attention (30 min.s) (Mohit)
- Aspect/opinion extraction (15 min.s) (Chelliah)
- RNNs (15 min.s) (Chelliah)

## **Morning break**

- Memory networks (25 min.s) (Asif)
- Code - Memory networks (30 min.s) (Mohit)
- Review Analyzer (15 min.s) (Mohit)
- RecursiveNN (10 min.s) (Chelliah)
- Convolutional Memory networks (20 min.s) (Chelliah)
- Cross-/Multi-lingual ABSA - (20 min.s) (Asif)
- Conclusion and Future Trends

# Sentiment Analysis

- Sentiment analysis aims to identify the orientation of opinion in a piece of text



# Why do we need Sentiment Analysis?

- *What others think* has always been an important piece of information
- *Overwhelming amount of information* on one topic: Manually reading or analysing all data is very inefficient
- *Biased/Fake* reviews
- An example
  - Mr. X needs to buy a phone. He was browsing amazon.in and found 1000 reviews for a particular phone.

## Scenario 1:

- Let there are **850 negative**, **100 positive** and **50 neutral** reviews
- Sentiment → **Negative**.

What if all the 100 positive reviews are at the top?

## Scenario 2:

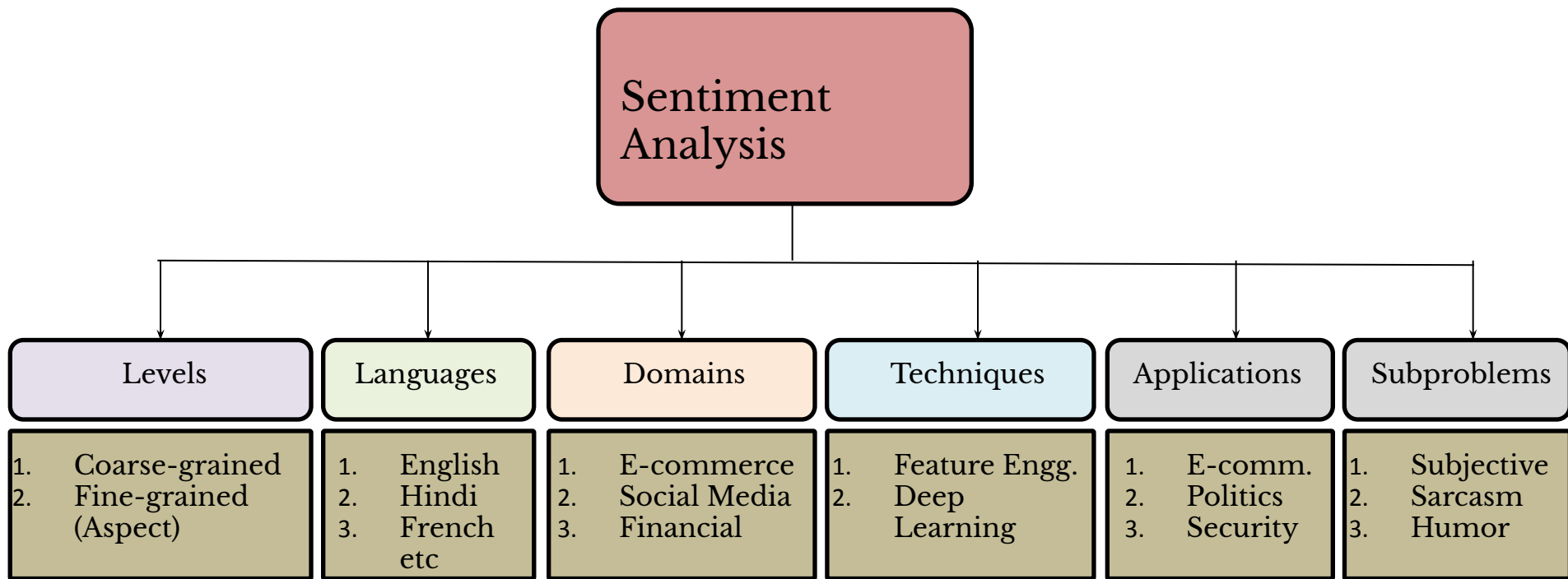
- Let there are **420 negative**, **480 positive** and **100 neutral** reviews.
- Sentiment → **Positive**

What if few of the reviews (e.g. 100) are fake? 4

# Challenges

- Similar lexical features but different sentiments
  - **This movie is not good**
  - **No movie can be better than this**
- Different style of writing but same sentiment
  - **It's an extremely useless phone**
  - **I have wasted my money on this phone**
  - **I could have bought Iphone instead of this**
- Product name, even, may appear in different forms
  - G-phone, Google-phone etc.
- Sentiment lexicons are not not sufficient for sentiment analysis
  - “*The food is very **cheap** here.*” vs “*The service is very **cheap** here.*”
- Reviews may not be genuine

# Sentiment Analysis: A broader view



# Sentiment Analysis: Subproblems

| Subproblems         | Text   | Remarks  |
|---------------------|--|--|
| Subjectivity        | <i>This movie is <b>awesome</b>.</i>   | Positive   |
|                     | <i>This movie is <b>pathetic</b>.</i>  | Negative   |
|                     | <i>This movie is <b>3-hours long</b>.</i>  | Neutral  |
| Thwarting           | <i><b>Impressive</b> story, <b>good</b> acting, however, it <b>didn't meet</b> my expectation.</i> | Small portion at the end dictates its sentiment.       |
| Sarcasm             | <i>This movie is <b>awesome to put you to sleep</b>.</i>   | Criticism in a humorous way.                           |
| Humble Bragging     | <i>My life is <b>miserable</b>, I have to <b>sign 300 autographs per day</b>.</i>                  | Draw attention to something of which someone is proud. |
| Discourse-based SA  | <i>This movie is a classic, <b>although</b>, I don't like 'sci-fi'.</i>                            | Sentiment is altered due to connectives.               |
| Sense-based SA      | <i>Shane Warne is a <b>deadly</b> spinner. (Positive)</i>  | Different sense leads to different sentiments.         |
|                     | <i>The campus has <b>deadly</b> snakes. (Negative)</i>   |  |
| Sentiment Intensity | <i>Movie was <b>ok</b>.</i>  | Weak positive sentiment.                               |
|                     | <i>Movie was <b>good</b>.</i>  | Mild positive sentiment.                               |
|                     | <i>Movie was <b>awesome</b>.</i>   | Strong positive sentiment.                             |

# Sentiment Analysis: Granularity

- Based on the granularity of analysis, we can categorize it as:
  - Coarse-grained Sentiment Analysis (Document-level or Sentence-level)
  - Fine-grained Sentiment Analysis (Phrase-level or Aspect-level)
- **Aspect Based Sentiment analysis (ABSA):** Sentiment towards an aspect (or opinion-target or feature)

Its **battery** is awesome but **camera** is very poor.

इसकी **बैटरी** शानदार है, लेकिन **कैमरा** बहुत ही खराब है।

(Isakee **baiTaree** shaanadaara hai, lekin **kaimaraa** bahut hee kharaab hai..)

**Positive** about the **battery** but **negative** about the **camera**



# Aspect Term Extraction

Given a set of sentences with pre-identified entities (e.g., restaurants), identify the aspect terms present in the sentence and return a list containing all the distinct aspect terms

“ I liked the **service** and the **staff** , but not the **food** ”

$\{ \textit{service}, \textit{staff}, \textit{food} \}$

“**Ambiance** and **music** funky, which I enjoy”

$\{ \textit{Ambiance}, \textit{music} \}$

“**Awesome** **form factor** and great **battery life**”

$\{ \textit{form factor}, \textit{battery life} \}$

# Polarity Identification

For a given set of aspect terms within a sentence, determine whether the polarity of each aspect term is *positive*, *negative*, *neutral* or *conflict* (i.e., both positive and negative)

I liked the *service* and the *staff* , but not the *food*

service: *Positive*, staff: *Positive*, food: *Negative*

I did add a *SSD drive* and *memory*

SSD drive: *Neutral*, memory: *Neutral*

The *RAM memory* is good but should have splurged for 8Mb instead of 4Mb

RAM memory: *Conflict*

# Aspect Based Sentiment Analysis: Few examples

*The **speed**, the **design**.. it is lightyears ahead of any PC I have ever owned.*

Speed, Design

Positive, Positive

***Tech support** would not fix the problem unless I bought your plan for \$150 plus.*

Tech support

Negative

*Certainly not the best **sushi** in New York, however, it is always fresh, and the **place** is very clean, sterile.*

Sushi, Place

Conflict, Positive

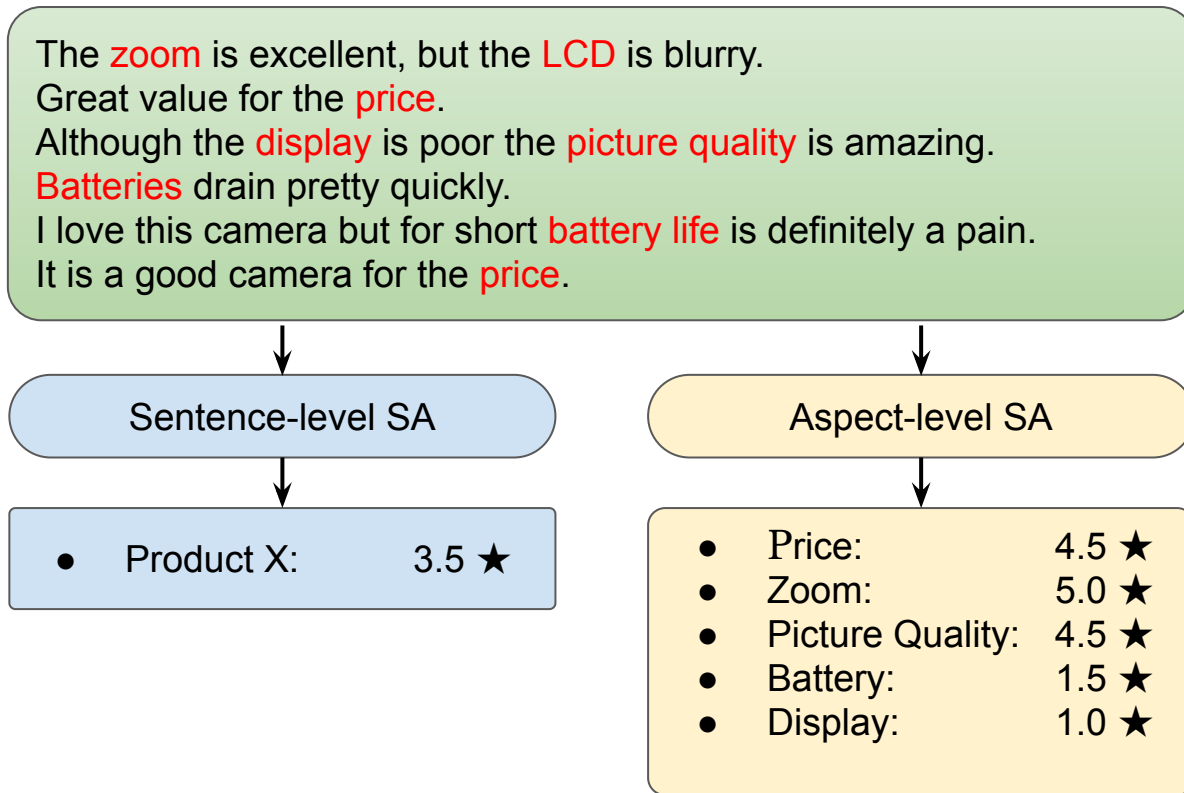
*It was very expensive for what you get.*

Implicit aspect (price), however, no textual presence implies no aspect term

*I enjoy having Apple products.*

-No Aspect term-

# Informed Decision: Coarse-grained vs Fine-grained SA



# ABSA: How attractive has been in recent time?

| Venue  | 2019 | 2018 | 2017 | 2016 | 2015 |
|--------|------|------|------|------|------|
| ACL    | 9    | 8    | 2    | 1    | 4    |
| EMNLP  | -    | 8    | 4    | 6    | 2    |
| COLING | -    | 6    | -    | 4    | -    |
| EACL   | -    | -    | 2    | -    | -    |
| NAACL  | 3    | 4    | -    | 1    | -    |
| AAAI   | 5    | 4    | 1    | 1    | 0    |
| IJCAI  | 6    | 4    | 1    | 1    | 1    |

Approaches to solve NLP problems

# Different Approaches for NLP

- Problems:
  - Part-of-Speech Tagging, Named Entity Recognition, Sentiment Analysis, Machine Translation
- Techniques:
  - Table driven system → Rule based system → Statistical system → Deep Learning system

# What technique should we use for NER?

## 1. Table driven system

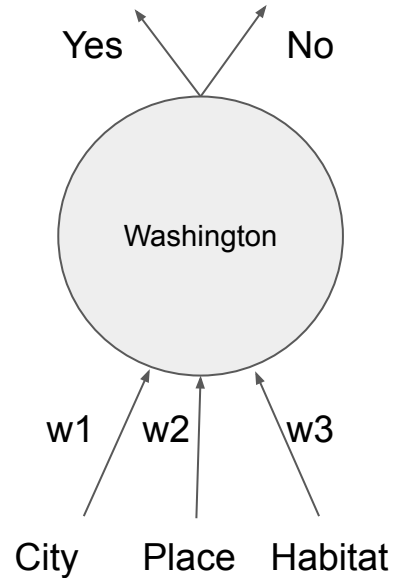
- Applicable only if there is no ambiguity
- Washington → Person or Location?

## 2. Rule based system

- Have to devise a number of rules (Kolmogorov complexity)
  - R1: if city in context then Location
  - R2: if capital in context then Location
  - ..
  - R1000: ...
- False positive for R1
  - Washington is a big city -- City in context, NER → Location
  - Washington was born in a big city -- City in context, NER → Person



# Rule based system: Representation as network



Inputs can be represented as one-hot vectors  
Weights are 1.

# What technique should we use?

## 3. Statistical

- $D = \operatorname{argmax} P(w|\text{context})$

## 4. Neural Network

- Key points:
  - Rules, Argmax and Neural Nets are interconvertible
  - Tables are lower bound.
  - Tables are much more universal.

# Is SA a table driven or machine learning approach?

- Whenever there is ambiguity table cannot help
- Types of SA
  - **Subjectivity:** *This movie is awesome.*
  - **Sarcasm:** *This movie is awesome to put you to sleep.*
  - **Thawarting:** *Impressive story, good acting, however, it didn't meet my expectation.*
  - **Humble bragging:** *My life is miserable, I have to sign 300 autographs per day.*
  - **Sense based SA:** *Shane Warne is a deadly spinner. v/s The campus has deadly snakes.*
  - **Discourse based SA:** *This movie is a classic, although, I don't like 'sci-fi'.*

# Traditional ML vs. DL pipeline

- Ngrams
- Presence or Absence of cue words
- Lexicons
- SVM
- Decision Tree



# RNN/LSTM based Aspect Sentiment Classification

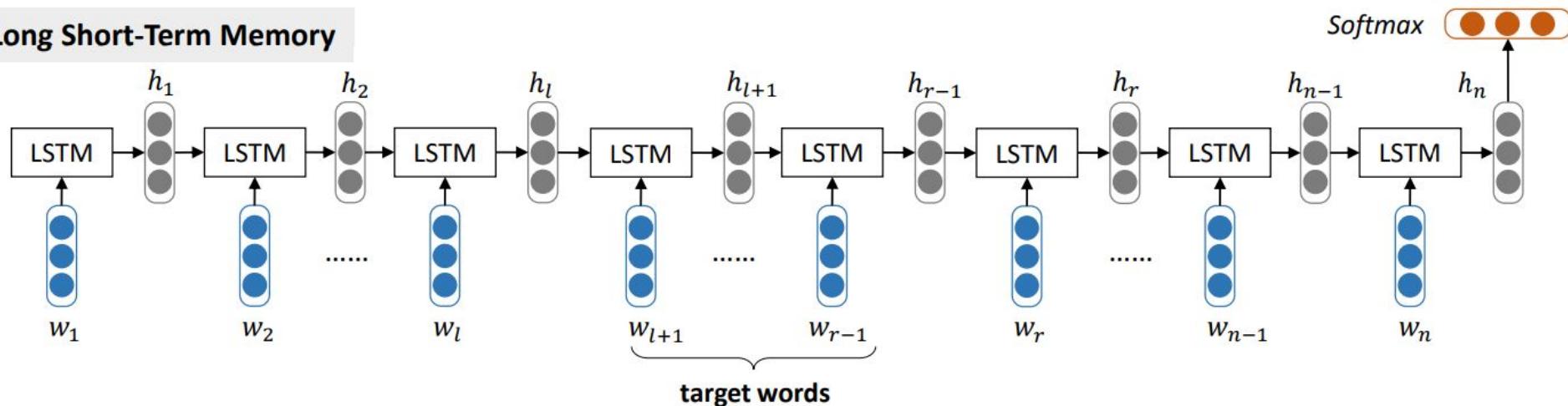
## Effective LSTMs for Target-Dependent Sentiment Classification [Tang et al. 2016]

- Long Short-Term Memory (LSTM)
  - Models the semantic representation of a sentence without considering the target word being evaluated
- Target-Dependent Long Short-Term Memory (TD-LSTM)
  - Extend LSTM by considering the target word
- Target-Connection Long Short-Term Memory (TC-LSTM)
  - Semantic relatedness of target with its context words are incorporated

## Simple LSTM [Tang et al. 2016]

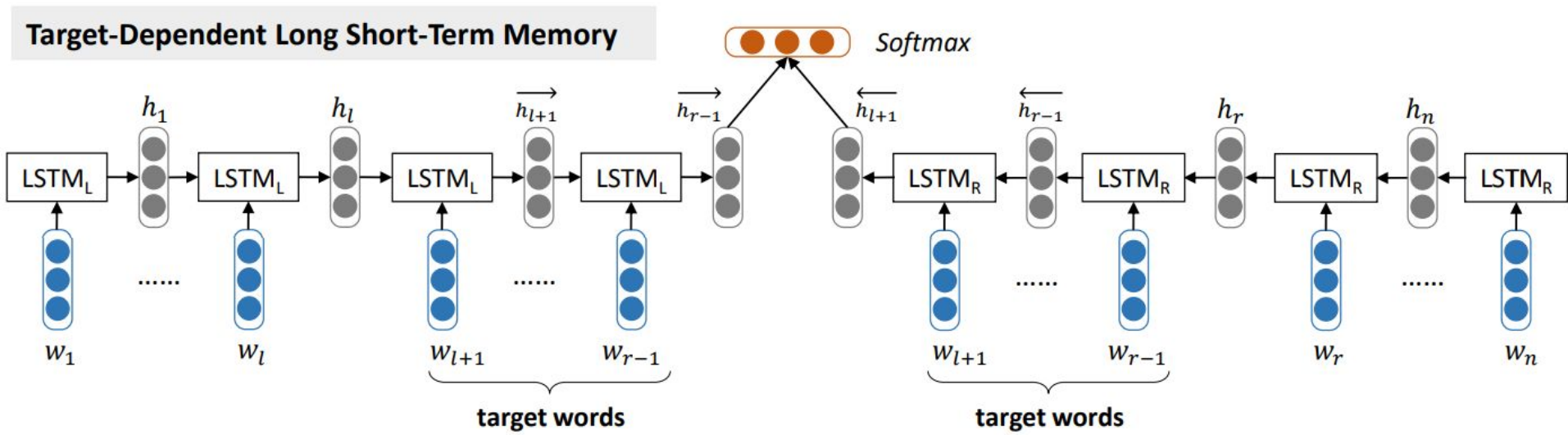
- Models the semantic representation of a sentence without considering the target word being evaluated
  - No discrimination between the following two instances
    - Its **battery** is awesome but camera is poor.
    - Its battery is awesome but **camera** is poor.

### Long Short-Term Memory



## Target-Dependent Long Short-Term Memory (TD-LSTM) [Tang et al. 2016]

- Considers the target word
  - Its **battery** is awesome but camera is poor.
    - $\text{LSTM}_L(\text{Its } \textbf{battery}) + \text{LSTM}_R(\text{battery is awesome but camera is poor.})$
  - Its battery is awesome but **camera** is poor.
    - $\text{LSTM}_L(\text{Its battery is awesome but } \textbf{camera}) + \text{LSTM}_R(\text{camera is poor.})$

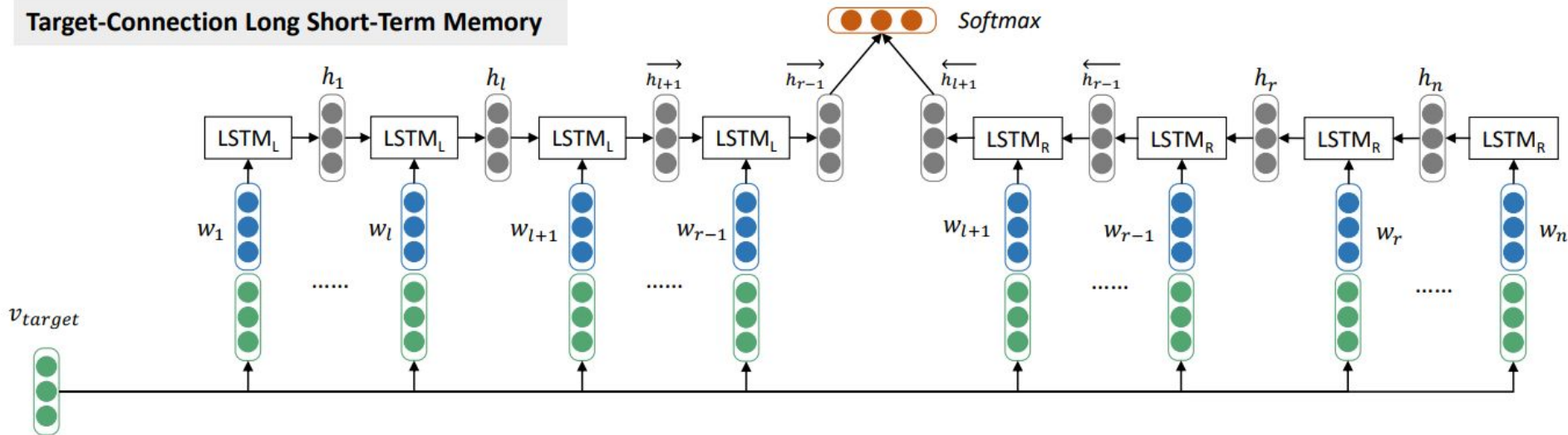




## Target-Connection Long Short-Term Memory (TC-LSTM) [Tang et al. 2016]

- Relationship between the word and the target is incorporated
  - Its **battery** is awesome but camera is poor.
    - $\text{LSTM}_{L_1}(\text{Its}, \text{battery}) \rightarrow \text{LSTM}_{L_2}(\text{battery}, \text{battery})$
    - $\text{LSTM}_{R_7}(\text{battery}, \text{battery}) \leftarrow \text{LSTM}_{R_6}(\text{is}, \text{battery}) \leftarrow \text{LSTM}_{R_5}(\text{awesome}, \text{battery}) \leftarrow \text{LSTM}_{R_4}(\text{but}, \text{battery}) \leftarrow \text{LSTM}_{R_3}(\text{camera}, \text{battery}) \leftarrow \text{LSTM}_{R_2}(\text{is}, \text{battery}) \leftarrow \text{LSTM}_{R_1}(\text{poor}, \text{battery})$

Target-Connection Long Short-Term Memory



# Experiments

- Dataset

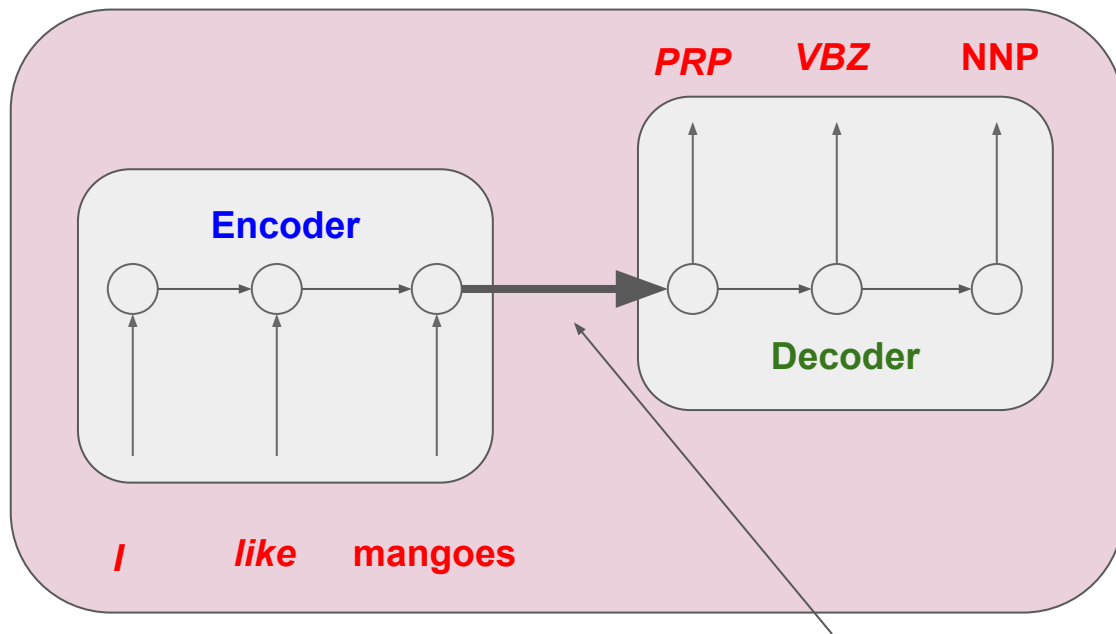
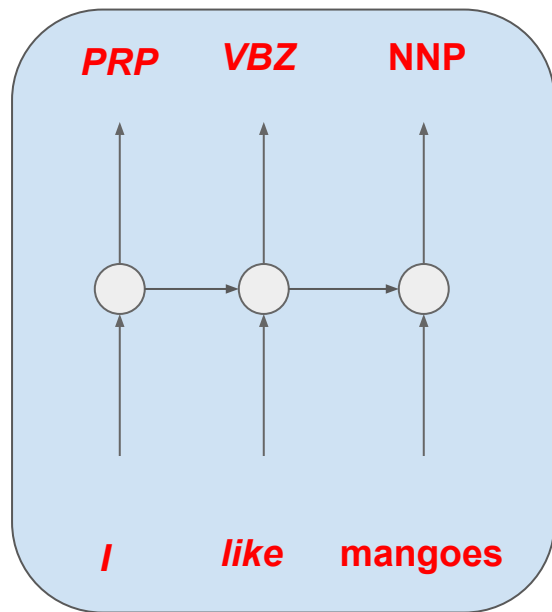
- Dong et al., 2014
  - Train: 6,248 sentences
  - Test: 692 sentences
  - Sentiment distribution: 25% → Positive, 25% → Negative, 50% → Neutral

| Method  | Accuracy | Macro-F1 |
|---------|----------|----------|
| LSTM    | 0.665    | 0.647    |
| TD-LSTM | 0.708    | 0.690    |
| TC-LSTM | 0.715    | 0.695    |

# Attention Mechanism

# Sequence labeling v/s Sequence transformation

- PoS Tagging



Sentence embeddings

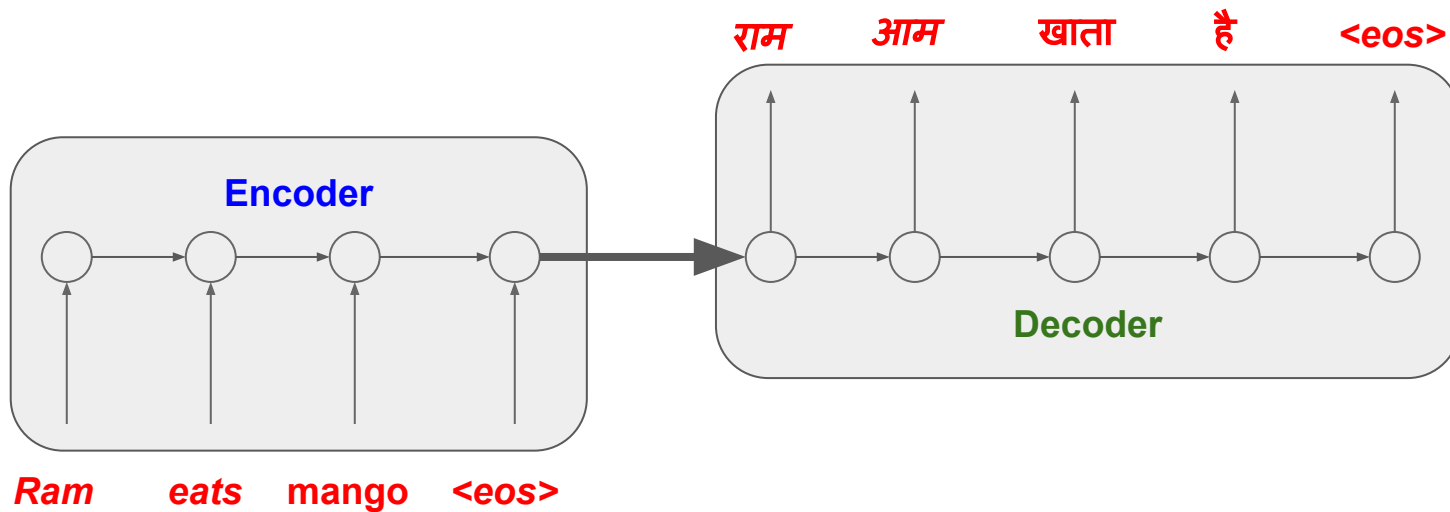
## Why is sequence transformation required?

- For many application length of I/p and O/p are not necessarily the same
  - E.g. *Machine Translation, Summarization, Question Answering* etc.
- For many applications length of O/p is not known
- Non-monotone mapping: Reordering of words
- PoS tagging, Named Entity Recognition etc. do not require these capabilities

# Encode-Decode paradigm

- English-Hindi Machine Translation

- Source sentence: 3 words
- Target sentence: 4 words
- Second word of the source sentence maps to 3rd & 4th words of the target sentence.
- Third word of the source sentence maps to 2nd word of the target sentence



## Problems with Encode-Decode paradigm

- Encoding transforms the entire sentence into a single vector
- Decoding process uses this sentence representation for predicting the output
  - Quality of prediction depends upon the quality of sentence embedding
- After few time steps decoding process may not properly use the sentence representation due to long-term dependency

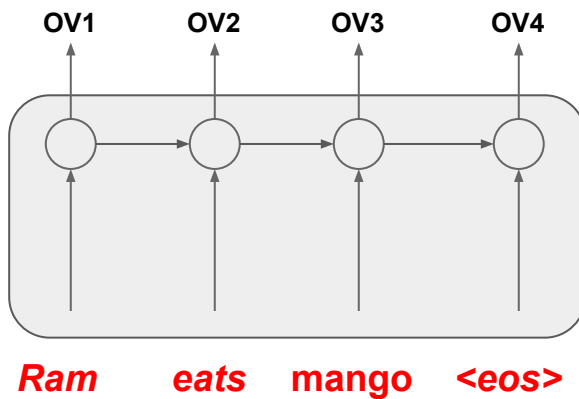
## Solutions

- To improve the quality of predictions we can
  - Improve the quality of sentence embeddings **'OR'**
  - Present the source sentence representation for prediction at each time step. **'OR'**
  - Present the RELEVANT source sentence representation for prediction at each time step.
    - *Encode - Attend - Decode* (Attention mechanism)



## Attention Mechanism

- Represent the source sentence by the set of **output vectors** from the encoder
- Each **output vector** (OV) at time  $t$  is a contextual representation of the input at time  $t$

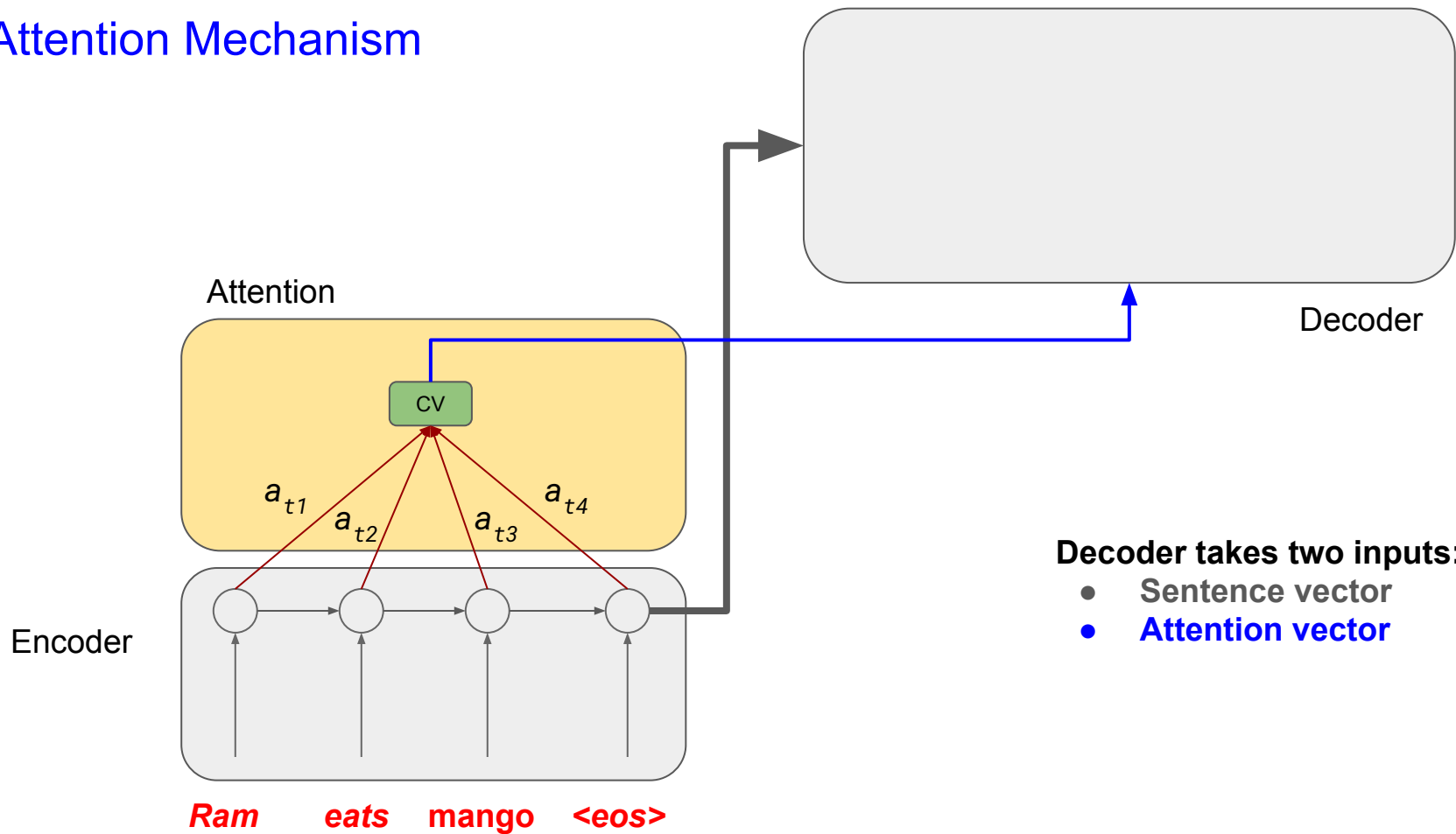


## Attention Mechanism

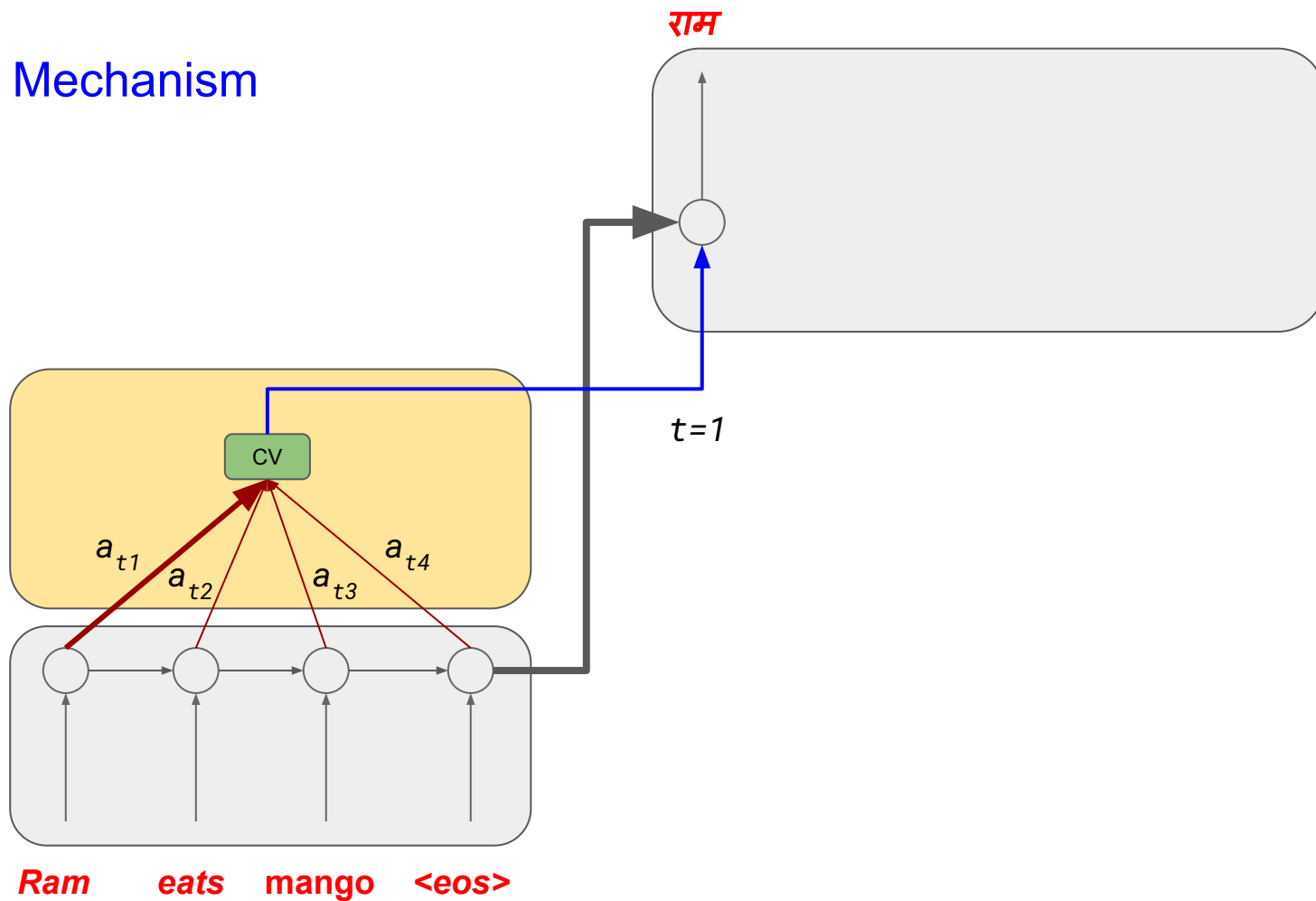
- Each of these output vectors (OVs) may not be equally relevant during decoding process at time  $t$ .
- Weighted average of the output vectors can resolve the relevancy.
  - Assign more weights to an output vector that needs more **attention** during decoding at time  $t$ .
- The weighted average **context vector (CV)** will be the input to decoder along with the sentence representation.
  - $CV_i = \sum a_{ij} \cdot OV_j$

where  $a_{ij}$  = weight of the  $j^{th}$  OV

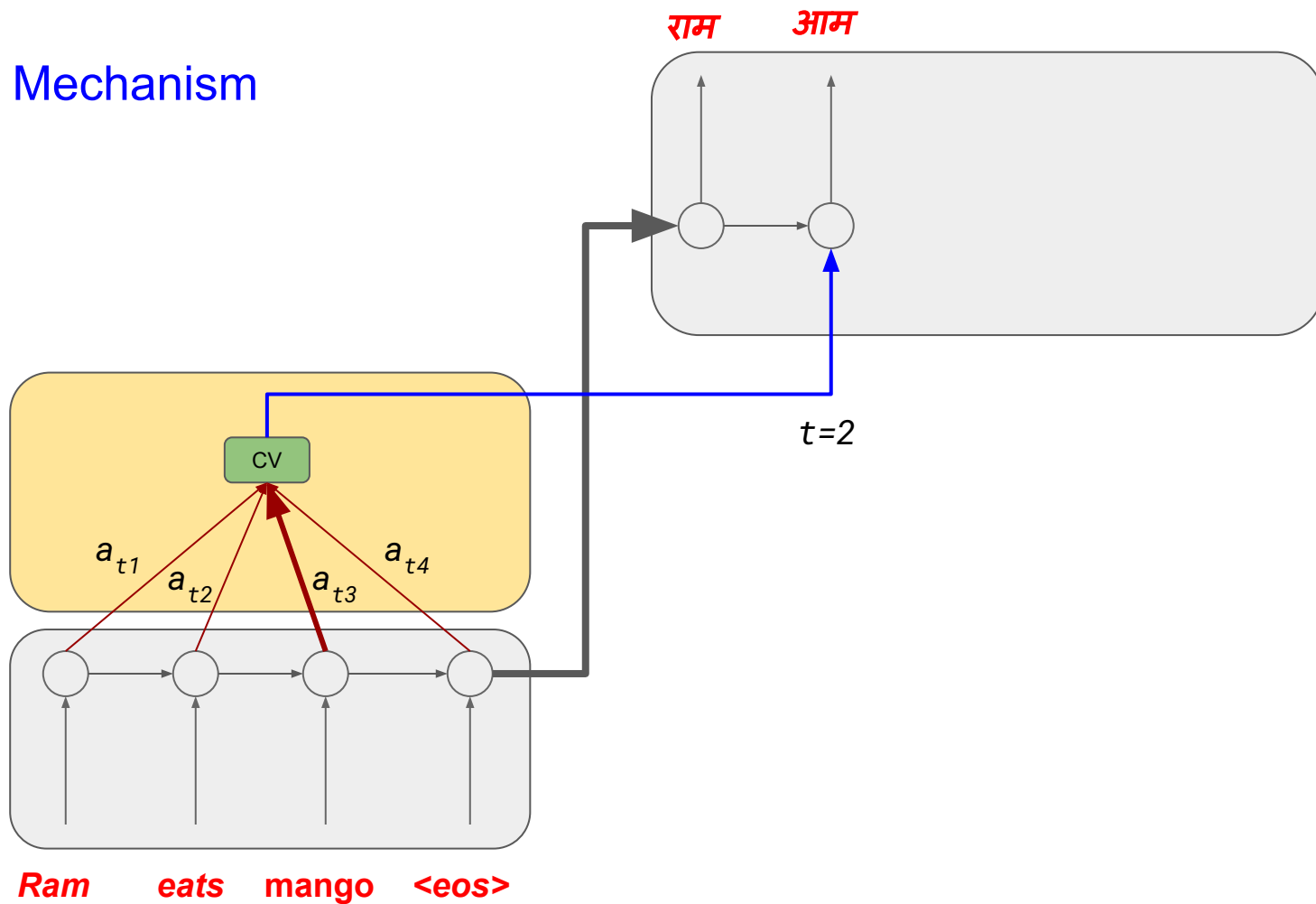
# Attention Mechanism



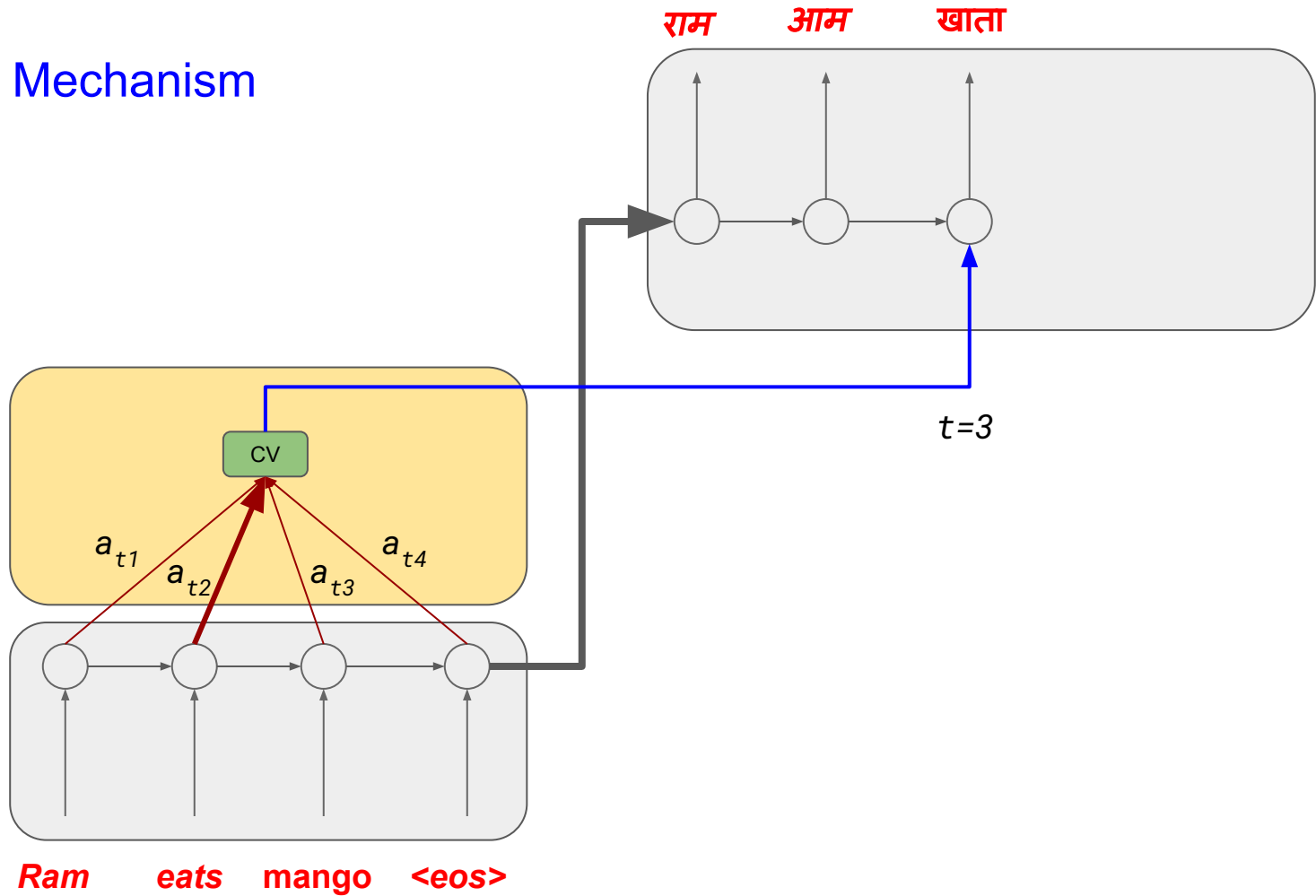
# Attention Mechanism



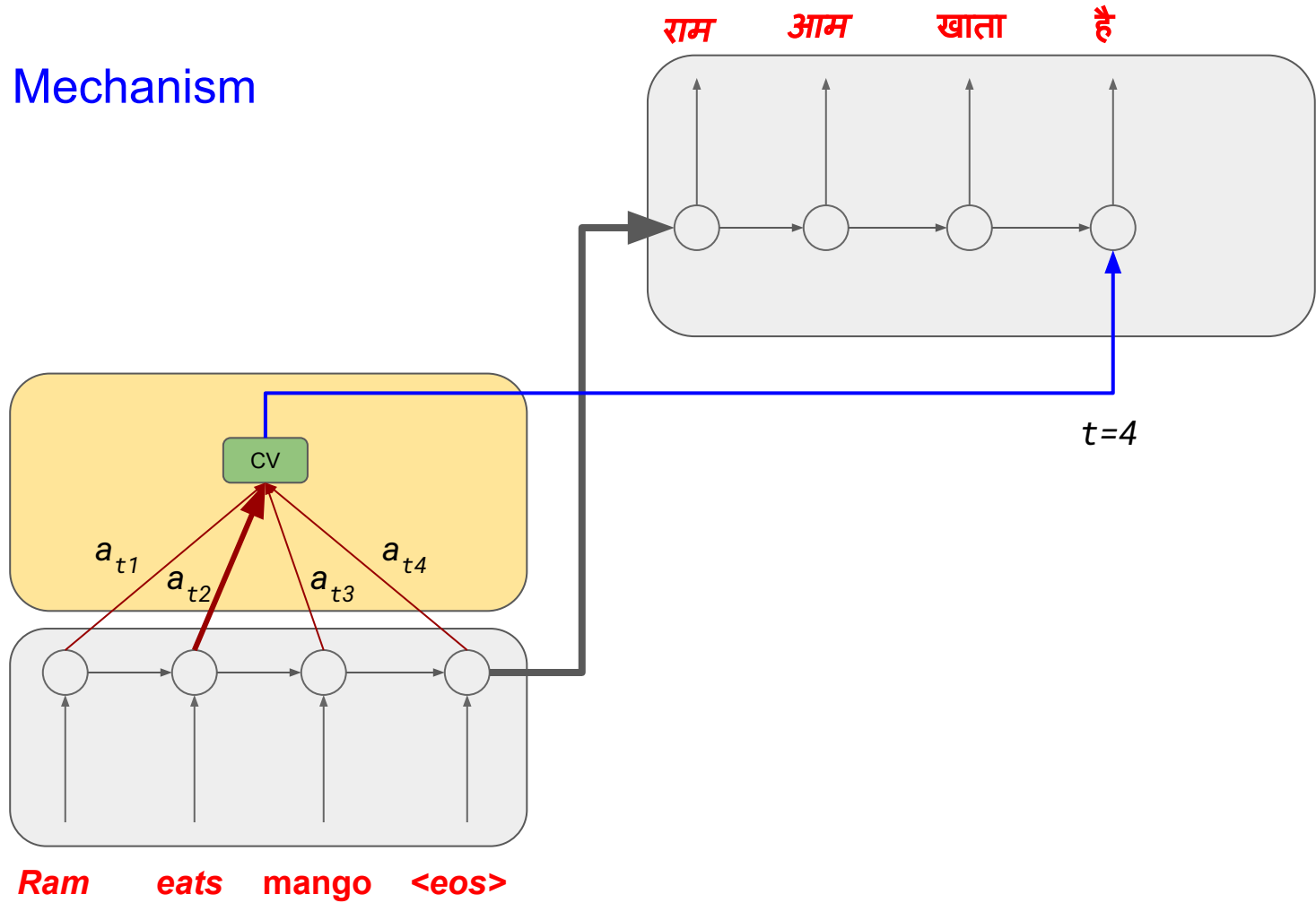
# Attention Mechanism



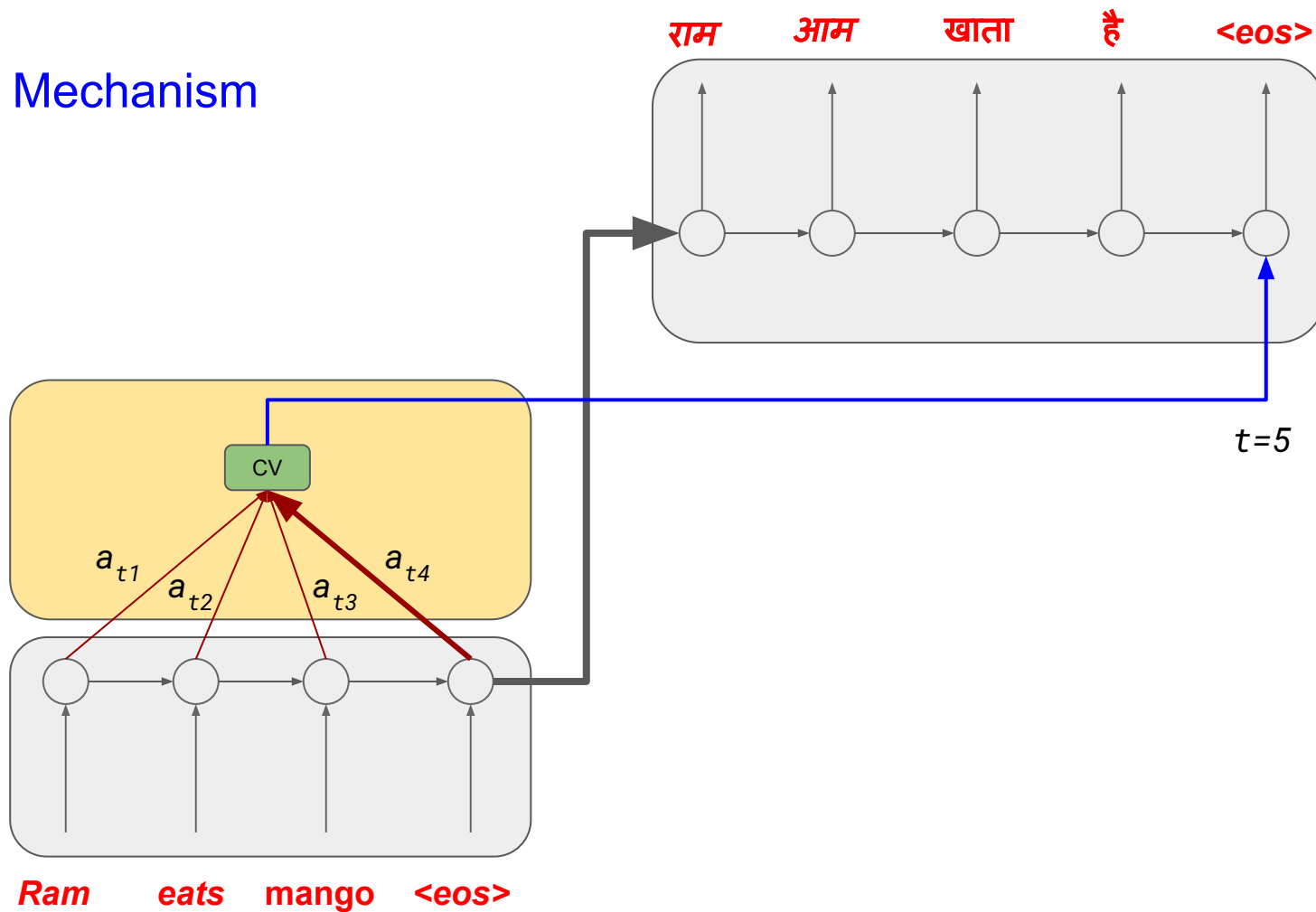
# Attention Mechanism



# Attention Mechanism



# Attention Mechanism





# Attention Mechanism for Classification

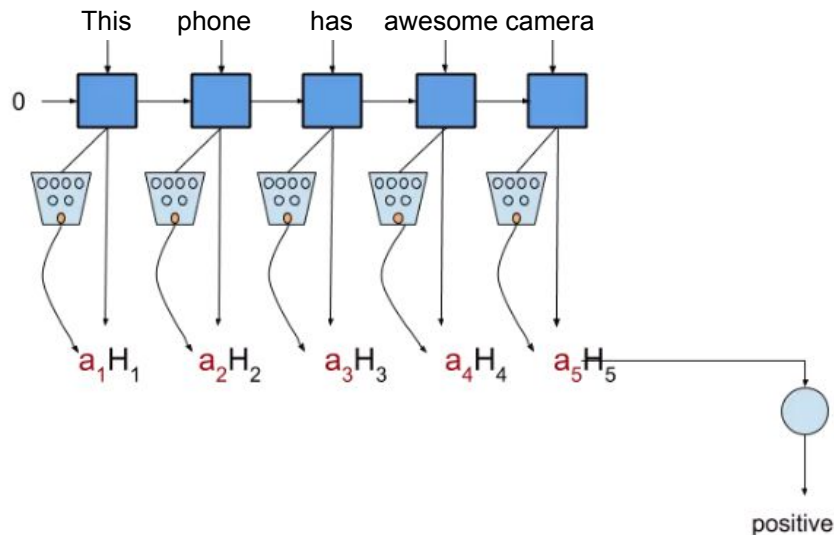
# Attention mechanism for classification

- Every word in a sentence is not equally important for any task
  - **Sentiment Classification:** *Adjectives* are more important than *prepositions* or *conjunctions*
    - This phone has awesome camera. → Word '*awesome*' is the most important word in the whole sentence considering the *positive sentiment*
- Why not weight each word in a sentence according to its importance?
- Attention mechanism is the solution
  - Compute attention weights ( $a_i$ ) by building a small fully-connected neural network on top of each encoded state
  - A single-unit final layer corresponds to the attention weight

$$y_i = \tanh(W \cdot H_i)$$

$$a_i = \exp(y_i) / \sum_j \exp(y_j)$$

$$h_i = a_i \cdot H_i$$



## Attention for Aspect Sentiment Classification

- Attend the important word considering the target
  - Its **battery** is awesome but **camera** is poor.
    - For target **battery**, awesome will have highest weight
    - For target **camera**, poor will have highest weight
- Through attention mechanism, the network can learn the association of *awesome* for **battery** and *poor* for **camera** in aspect sentiment classification.

## Attention-based LSTM for Aspect-level Sentiment Classification [Wang et al. 2016]

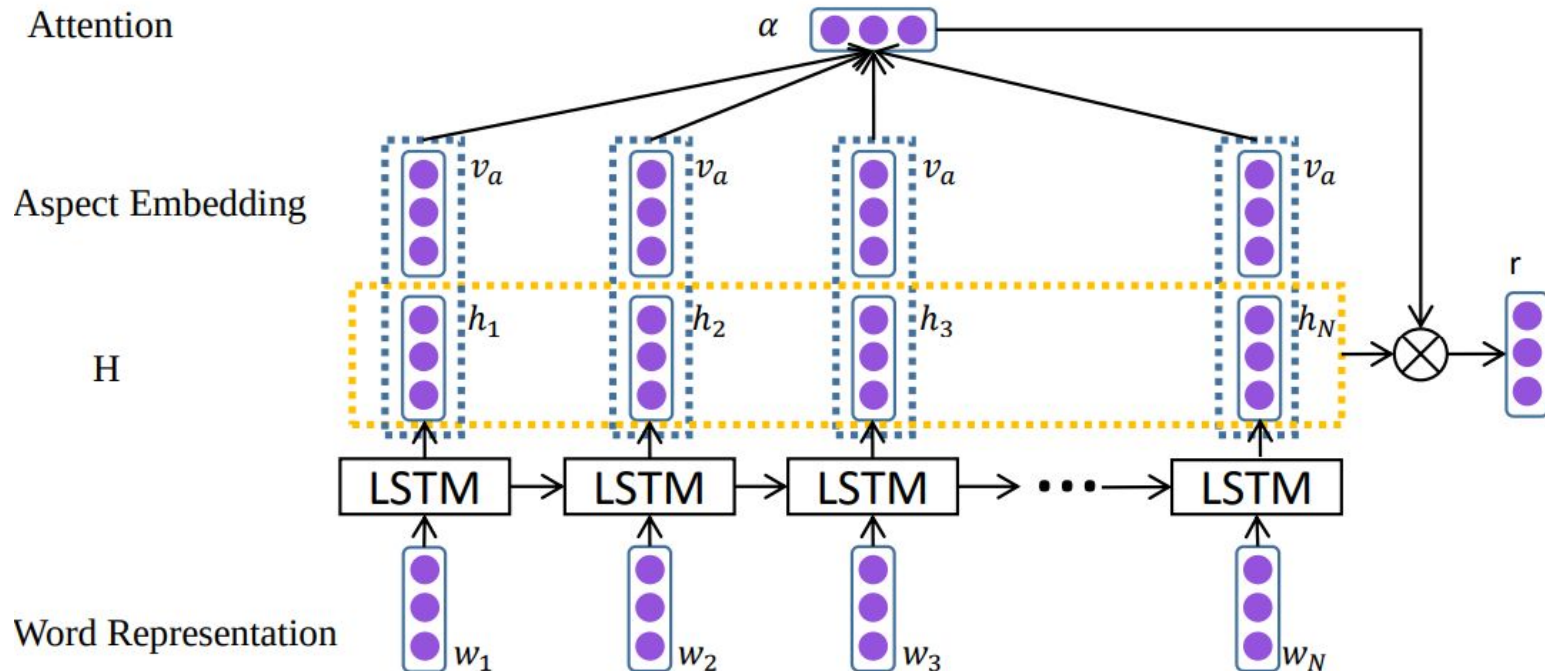
- Incorporation of only target information is not sufficient
- Application of attention mechanism can extract the association of important word for an aspect
- Two architectures
  - Attention-based LSTM (AT-LSTM)
    - Relationship between the word and the target is incorporated at the attention layer ONLY
  - Attention-based LSTM with Aspect Embedding (ATAE-LSTM)
    - Relationship between the word and the target is incorporated at the input and attention layer BOTH

Aspect  $a = \text{battery life}$

Aspect Embedding  $v_a = (\text{emb}_{\text{battery}} + \text{emb}_{\text{life}}) / 2$

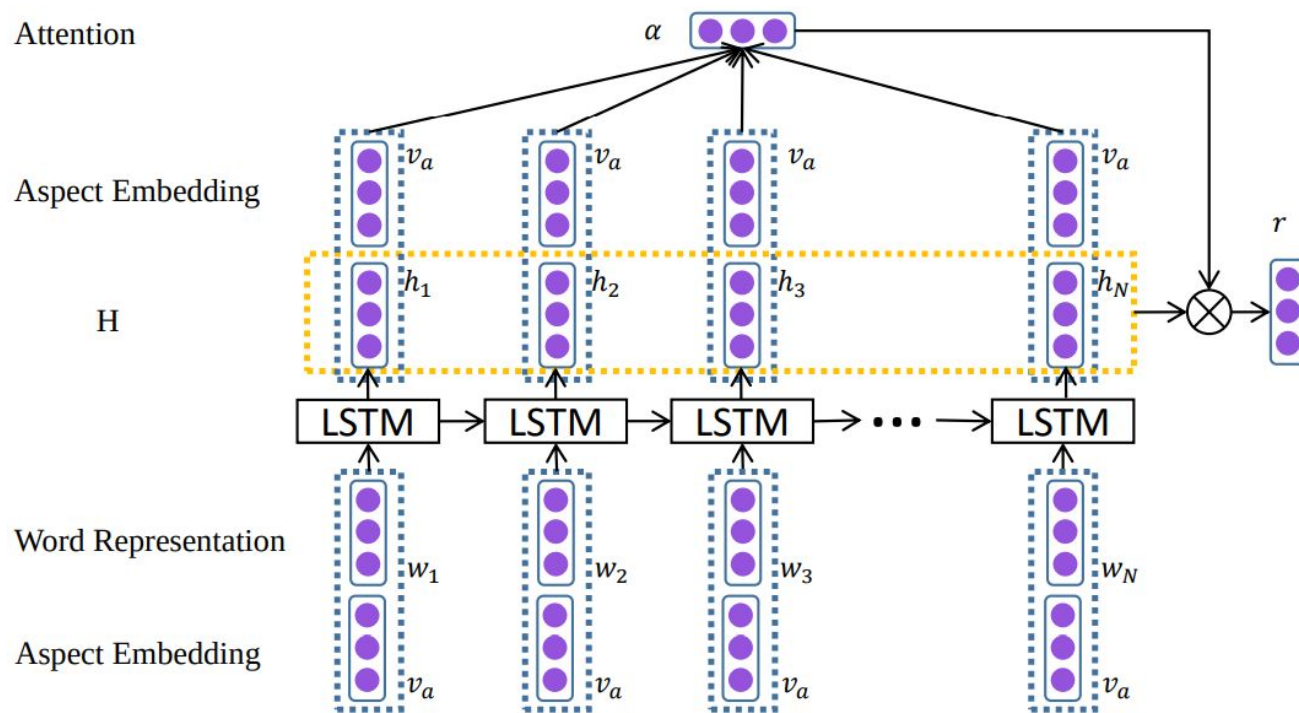
## Attention-based LSTM (AT-LSTM)

- Relationship between the word and the target is incorporated at the attention layer



## Attention-based LSTM with Aspect Embedding (ATAE-LSTM)

- Relationship between the word and the target is incorporated at the input layer and the attention layer



## Datasets

- SemEval-2014 [Pontiki et al., 2014]

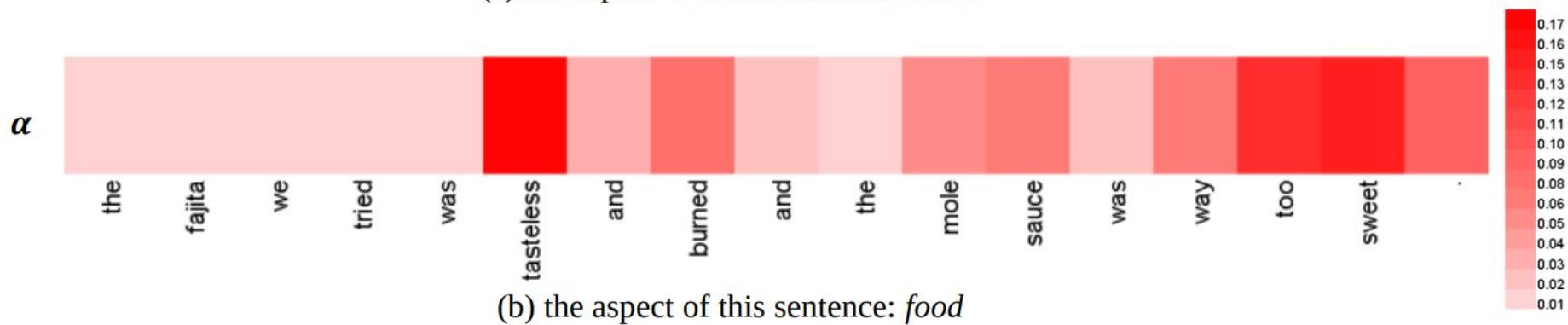
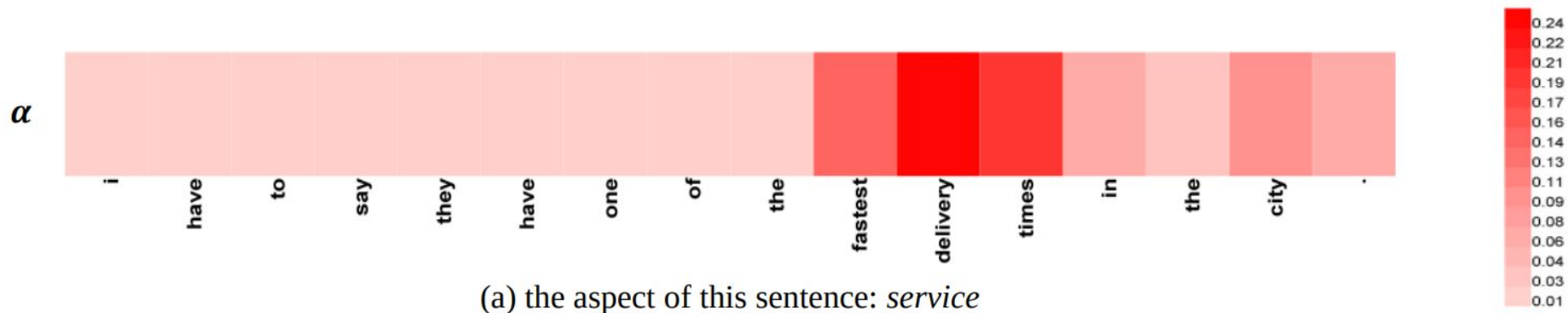
| Aspect   | Positive |      | Negative |      | Neutral |      |
|----------|----------|------|----------|------|---------|------|
|          | Train    | Test | Train    | Test | Train   | Test |
| Food     | 867      | 302  | 209      | 69   | 90      | 31   |
| Price    | 179      | 51   | 115      | 28   | 10      | 1    |
| Service  | 324      | 101  | 218      | 63   | 20      | 3    |
| Ambience | 263      | 76   | 98       | 21   | 23      | 8    |
| Misc     | 546      | 127  | 199      | 41   | 357     | 51   |
| Total    | 2179     | 657  | 839      | 222  | 500     | 94   |

## Experimental Results

| Method    | Pos/Neg/Neu | Pos/Neg |
|-----------|-------------|---------|
| LSTM      | 82.0        | 88.3    |
| TD-LSTM   | 82.6        | 89.1    |
| TC-LSTM   | 81.9        | 89.2    |
| AT-LSTM   | 83.1        | 89.6    |
| ATAE-LSTM | 84.0        | 89.9    |



## Attention weights: Heatmaps



## Aspect-level Sentiment Classification with HEAT (HiErarchical ATtention) Network [Cheng et al. 2017]

- Introduced HiErarchical ATtention (HEAT) network
  - Aspect attention (*with respect to the aspect category*)
  - Sentiment attention (*with respect to the aspect based attention*)
- Aspect attention
  - pays attention to the aspect information, i.e., aspect terms, under the direction of the target aspect
- Sentiment attention
  - aims to capture the sentiment feature of the text under the direction of the target aspect and the extracted aspect information

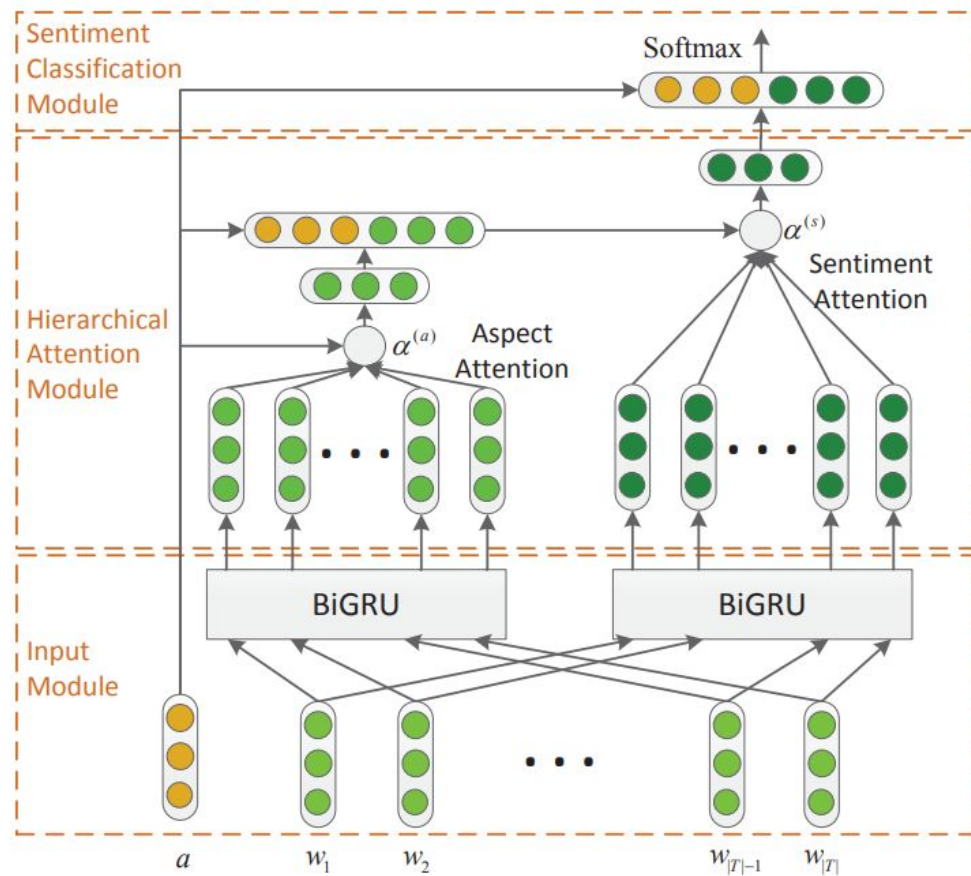
# A motivating example

The tastes are great, but the service is dreadful

- Both sentiment-bearing words, *great* and *dreadful* can be used for both aspects, *food* and *service*
- Given aspect *food*, model can attend to both *great* and *dreadful*- *Confusing !*
- Remedy
  - Leveraging aspect term to bridge the gap
  - Given aspect *food*, much easier to find aspect term *tastes* than to discriminate which sentiment word is corresponding to the aspect (*through aspect attention*)
  - Under the guidance of aspect term *tastes*, we can easily choose the sentiment word *great* and decide the sentiment polarity on the aspect

# HEAT for Aspect Sentiment Classification

- Aspect attention aims to pay attention to the aspect information, i.e., aspect terms (**taste**), under the direction of the target aspect (**food**)
- Sentiment attention aims to capture the sentiment feature of the text (**great**) under the direction of the target aspect (**food**) and the extracted aspect information (**tastes**)



# Experiments

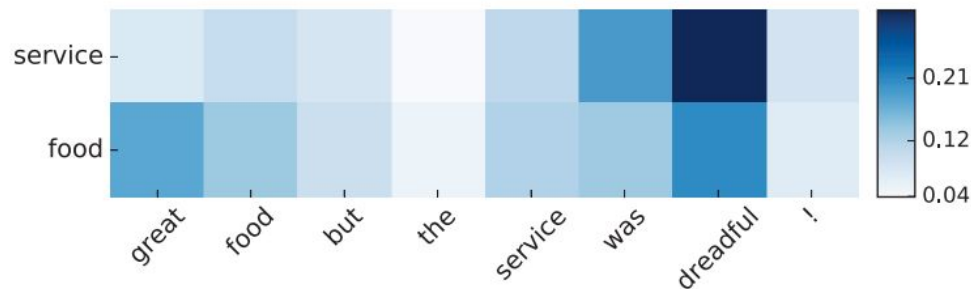
- Dataset

- SemEval-2014 [Pontiki et al., 2014] → Restaurant
- SemEval-2015 [Pontiki et al., 2015] → Restaurant and Laptop
- SemEval-2016 [Pontiki et al., 2016] → Restaurant

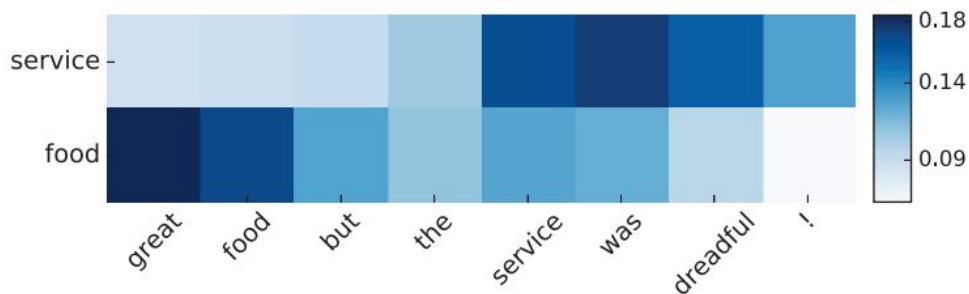
HEAT: Standard Attention (Softmax)  
HEATB: Bernoulli Attention (Sigmoid)

| Method      | Restaurant 14 |             | Restaurant 15 |             | Restaurant 16 |             | Laptop 15 |             |
|-------------|---------------|-------------|---------------|-------------|---------------|-------------|-----------|-------------|
|             | Pos/Neg       | Pos/Neg/Neu | Pos/Neg       | Pos/Neg/Neu | Pos/Neg       | Pos/Neg/Neu | Pos/Neg   | Pos/Neg/Neu |
| AT-LSTM     | 89.6          | 83.1        | 81.0          | 77.2        | 87.6          | 83.0        | 86.3      | 82.1        |
| ATAE-LSTM   | 89.9          | 84.0        | 80.9          | 77.4        | 87.2          | 82.7        | 85.8      | 82.3        |
| AT-BiGru    | 90.4          | 84.3        | 82.8          | 79.2        | 90.4          | 86.7        | 87.0      | 84.3        |
| HEAT-GRU    | 89.6          | 84.3        | 81.2          | 79.1        | 89.7          | 85.5        | 87.8      | 84.5        |
| HEATB-GRU   | 89.4          | 84.0        | 81.8          | 79.6        | 89.2          | 85.4        | 87.3      | 84.2        |
| HEAT-BiGRU  | 91.3          | 85.1        | 83.0          | 80.1        | 90.8          | 87.1        | 87.9      | 84.9        |
| HEATB-BiGRU | 91.1          | 84.9        | 83.4          | 80.5        | 91.1          | 87.5        | 88.0      | 85.1        |

## Attention Analysis



(a) Result of AT-BiGRU.



(b) Result of HEATB-BiGRU.

1. AT-BiGRU gets confused to locate sentiment word for aspect food in Figure 4(a)

Given aspect food, both “*great*” and “*dreadful*” obtain high scores

2. In Figure 4(b) HEATB-BiGRU solves the problem well

Expression “*service was dreadful!*” gets higher scores than other words given aspect *service*

Expression “*great food*” achieves the top scores given aspect *food*

## Interactive Attention Networks for Aspect-Level Sentiment Classification [Ma et al. 2017]

- Previous approaches incorporated the target information (i.e. *aspect*) for modelling the target-specific contexts
  - Generated target-specific representations
- Studies ignored the separate modeling of target with respect to context
- BUT, coordination of targets and contexts could be useful
  - Example, “*The picture quality is clear-cut but the battery life is too short*”

When *short* is collocated with *battery life*, sentiment class is *negative*

- BUT, for *Short fat noodle spoon, relatively deep some curva*

When *short* is collocated with *spoon*, sentiment tends to be *neutral*

## Interactive Attention Networks for Aspect-Level Sentiment Classification [Ma et al. 2017]

- Now, the issue

### *How to simultaneously model the target and context precisely?*

- First, target and context can determine representations of each other

For example, when we see the target “**picture quality**”, context word “**clear-cut**” is naturally associated with the target and the vice-versa

*We argue that targets and contexts can be modeled separately but learned from their interaction*

- Second, different constituents of a target aspect and context offer different information

For example, it is easy to know that “**picture**” plays a more important role in the representation of the target “**picture quality**” (described by **clear-cut**)

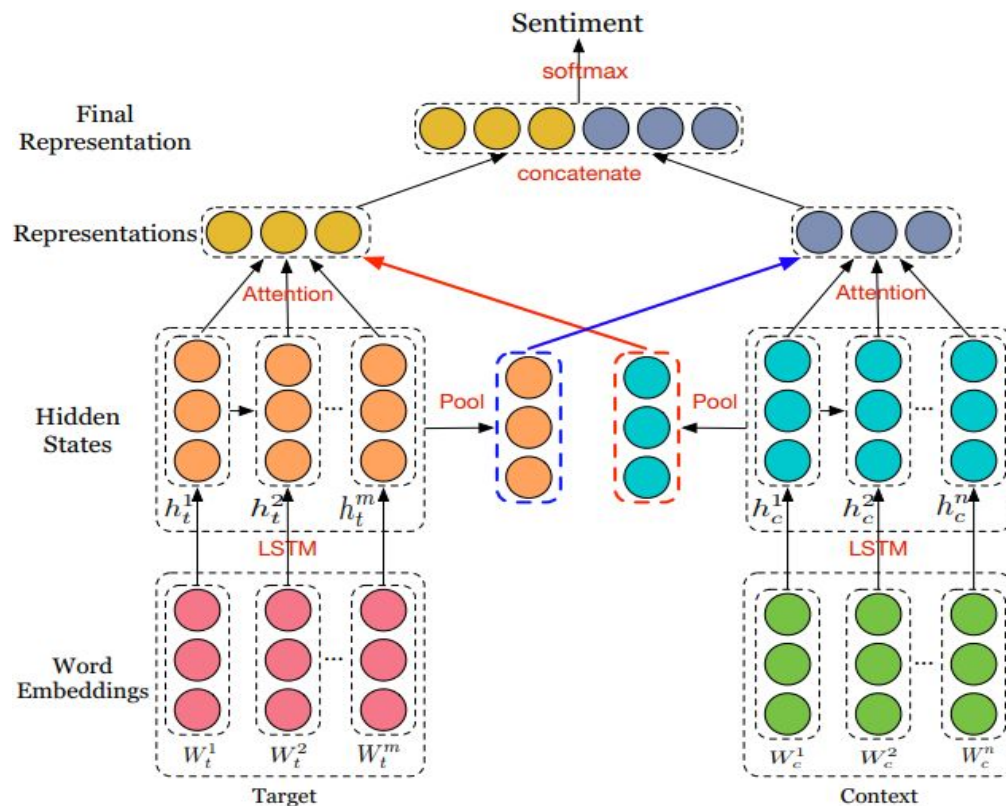


## Interactive Attention Networks for Aspect-Level Sentiment Classification [Ma et al. 2017]

- Both *targets* and *contexts* deserve special treatment and need to be learned their own representations via *interactive learning*
- Why interactive?
  - Interactively learn attentions in the contexts and targets, and generate the representations for targets and contexts separately
- **Steps of IAN**
  - Utilizes the attention mechanism associated with a target to get important information from the context and compute context representation for sentiment classification
  - Makes use of the interactive information from context to supervise the modeling of the target which is helpful to judging sentiment
  - Finally, with both target representation and context representation concatenated, IAN predicts the sentiment polarity for the target within its context

# Interactive Attention Networks (IAN)

- IAN learns the attentions for the contexts and targets separately
  - Generates the separate representations for targets and contexts via interaction with each other



## Experimental Results

- Dataset: SemEval-2014 [Pontiki et al., 2014]
  - Restaurant and Laptop

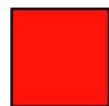
| Method    | Restaurant   | Laptop       |
|-----------|--------------|--------------|
| Majority  | 0.535        | 0.650        |
| LSTM      | 0.743        | 0.665        |
| TD-LSTM   | 0.756        | 0.681        |
| AE-LSTM   | 0.762        | 0.689        |
| ATAE-LSTM | 0.772        | 0.687        |
| IAN       | <b>0.786</b> | <b>0.721</b> |

## Attention weights: Heatmap

“The ***fish*** is fresh but the ***variety of fish*** is nothing out of ordinary.

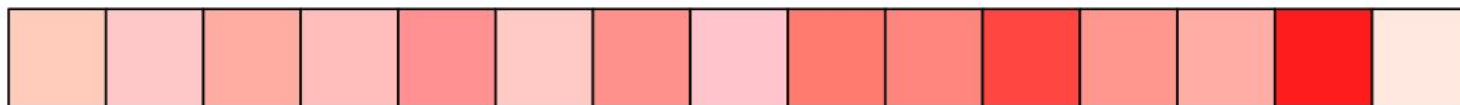


variety of fish

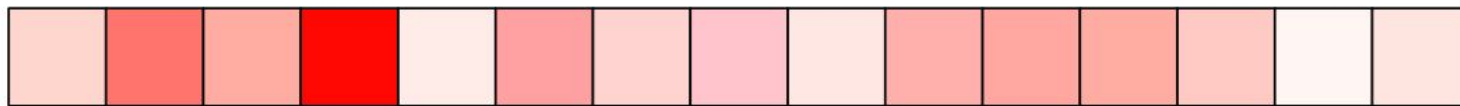


fish

(a) weight for target



the fish is fresh but the variety of fish is nothing out of ordinary .



the fish is fresh but the variety of fish is nothing out of ordinary .

(b) weight for context

# Effective Attention Modeling for Aspect-Level Sentiment Classification [He et al. 2018]

- Improved the *effectiveness of attention mechanism* to capture the *importance of each context* word towards a *target* by modeling their semantic associations
  - Proposed a method for ***target representation*** that better captures the semantic meaning of the ***opinion target***
  - Introduced an ***attention model*** that incorporates ***syntactic information*** obtained from a ***dependency parser***

## Target Representation

- While computing attention, simple averaging may not capture the real semantics of the target well
  - E.g., “*hot dog*” → Averaging of vectors may not represent it closer to the cluster of food items
- Represent the target as a weighted summation of aspect embeddings
- For aspect embedding matrix  $T \in R^{K \times d}$ , the target representation is computed as follows

$$\mathbf{t}_s = \mathbf{T}^\top \cdot \mathbf{q}_t$$

$$\mathbf{q}_t = \text{softmax}(\mathbf{W}_t \cdot \mathbf{c}_s + \mathbf{b}_t)$$

$$\mathbf{c}_s = \text{Average}\left(\frac{1}{m} \sum_{i=1}^m \mathbf{e}_{a_i}, \frac{1}{n} \sum_{j=1}^n \mathbf{e}_{w_j}\right)$$

where  $K$  is number of predefined aspects (e.g., food, price, service, ambience and misc),  $m$  is the length of target,  $n$  is the length of sentence and  $\mathbf{e}$  stands for embedding

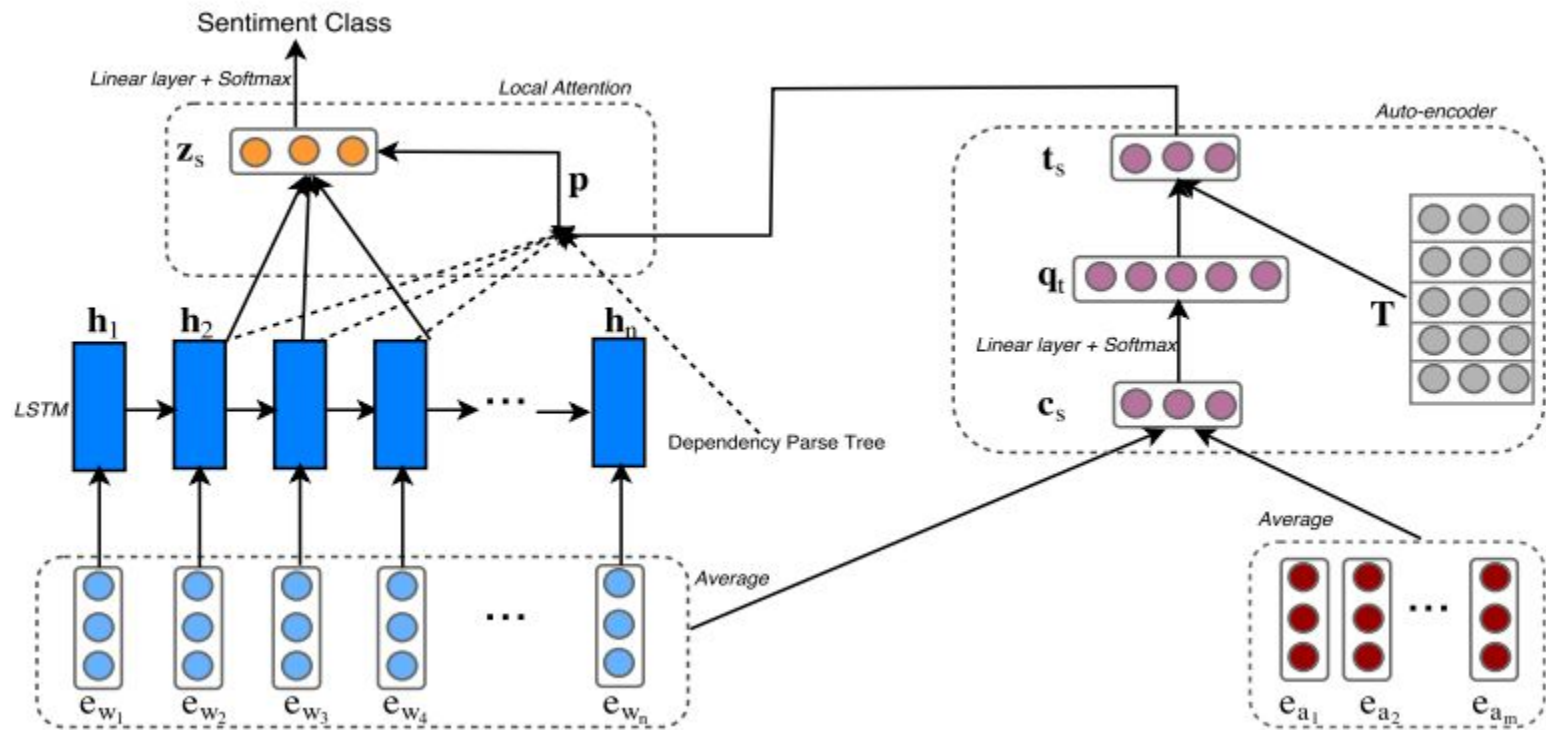
## Syntactic Information

- *Opinion words* that are closer to the *target* in the *dependency tree* are more relevant for determining its sentiment
- Attention model selectively attends to a small window of context words based on their location
- For a context window  $ws$ ,

$$p_i = \frac{d_i}{\sum_j d_j}$$
$$d_i = \begin{cases} \frac{1}{2^{(l_i-1)}} \cdot \exp(f_{score}(\mathbf{h}_i, \mathbf{t}_s)), & \text{if } l_i \in [1, ws] \\ 0, & \text{otherwise} \end{cases}$$
$$f_{score}(\mathbf{h}_i, \mathbf{t}_s) = \tanh(\mathbf{h}_i^T \cdot \mathbf{W}_a \cdot \mathbf{t}_s)$$

where  $\mathbf{t}_s$  is the target representation,  $l_i$  is the distance from the target in the dependency tree

# Architecture





# Experiments

- Dataset

- SemEval-2014 [Pontiki et al., 2014] → Restaurant and Laptop
- SemEval-2015 [Pontiki et al., 2015] → Restaurant
- SemEval-2016 [Pontiki et al., 2016] → Restaurant

| Method              | Restaurant 14 |              | Laptop 14    |              | Restaurant 15 |              | Restaurant 16 |              |
|---------------------|---------------|--------------|--------------|--------------|---------------|--------------|---------------|--------------|
|                     | Accuracy      | Macro-F1     | Accuracy     | Macro-F1     | Accuracy      | Macro-F1     | Accuracy      | Macro-F1     |
| SVM                 | 80.16         | NA           | 70.49        | NA           | NA            | NA           | NA            | NA           |
| LSTM                | 75.23         | 64.21        | 66.79        | 64.02        | 75.28         | 54.1         | 81.94         | 58.11        |
| LSTM+Attn           | 76.83         | 66.48        | 68.07        | 65.27        | 77.38         | 60.52        | 82.73         | 59.12        |
| TDLSTM              | 75.37         | 64.51        | 68.25        | 65.96        | 76.39         | 58.7         | 82.16         | 54.21        |
| TDLSTM+Attn         | 75.66         | 65.23        | 67.82        | 64.37        | 77.1          | 59.46        | 83.11         | 57.53        |
| ATAE-LSTM           | 78.6          | 67.02        | 68.88        | 65.93        | 78.48         | 62.84        | 83.77         | 61.71        |
| MemNet              | 76.87         | 66.4         | 68.91        | 63.95        | 77.89         | 59.52        | 83.04         | 57.91        |
| LSTM+Attn+TarRep    | 78.95         | 68.67        | 70.69        | 66.59        | 80.05         | <b>68.73</b> | 84.24         | <b>68.62</b> |
| LSTM+SynAttn        | 80.45         | 71.26        | <b>72.57</b> | 69.13        | 80.28         | 65.46        | 83.39         | 66.83        |
| LSTM+SynAttn+TarRep | <b>80.63</b>  | <b>71.32</b> | 71.94        | <b>69.23</b> | <b>81.67</b>  | 66.05        | <b>84.61</b>  | 67.45        |

## Hierarchical Attention based Position-aware Network for Aspect-level Sentiment Analysis [Li et al. 2018]

- Introduces **position embeddings** to learn the **position-aware representations** of sentences and generate the **target-specific representations** of contextual words
- Position of a target aspect in a sentence provides useful evidence
  - *“I bought a mobile phone, its camera is wonderful but the battery life is a bit short”*
  - In context window approach: For *“battery life”*, both *“wonderful”* and *“short”* are likely to be considered as its adjunct word
  - If we encode the position information into the representation of each word effectively, we would have more confidence in concluding that the **“short”** is the adjunct word of **“battery life”** and predict the **sentiment** as **negative**
- Encode the position information into the representation of each word effectively

# Hierarchical Attention Based Position-aware Network (HAPN)

- Position embeddings of word  $w_i$

where

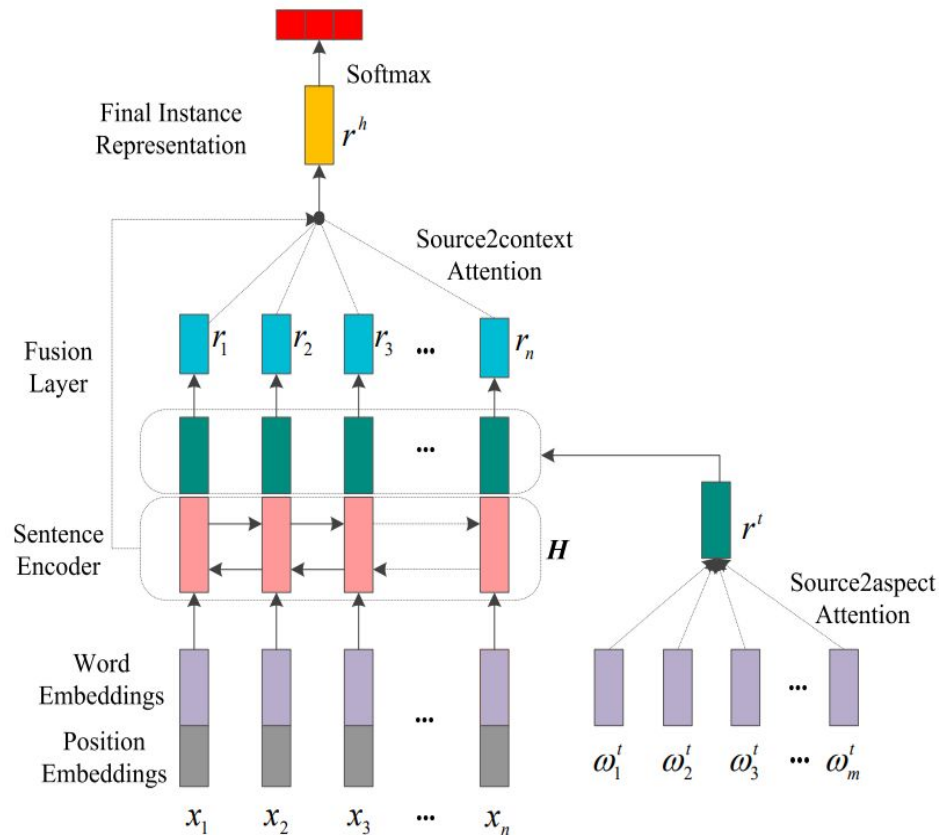
$$\begin{cases} i - k & i < k \\ i - k - m & n \geq i > k + m \\ 0 & k + m \geq i \geq k \end{cases}$$

$k$ : Index of first word of target

$m$ : Length of the target

$n$ : Length of the sentence

- Source2aspect** Attention
  - capture the most important clues in the target words
- Source2context** Attention
  - capture the most indicative sentiment words in the context
  - generates weighted-sum embedding for sentence representation



# Hierarchical Attention: More details

- Source2Aspect attention
  - Similar to self-attention
  - Generates the representation of aspect
  - Subsequently, *aspect-specific representation of each word = aspect representation + encoded position-aware representation*
  - Position-aware encoding: corresponds to the output of Bi-GRU that has input as *position embedding + word embedding*
- Source2Context attention
  - Captures the most indicative sentiment words in the context
  - Generates the weighted sum embeddings as the final sentence representation

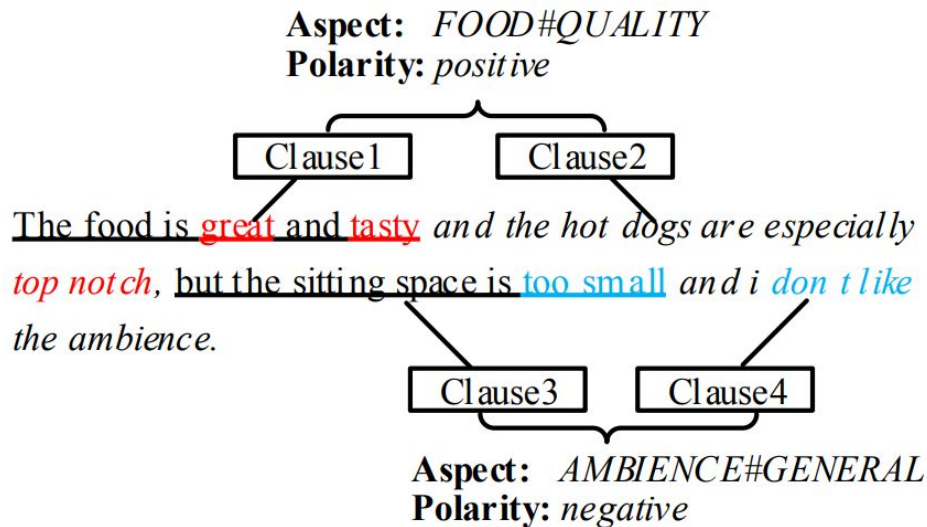
# Experiments

- Dataset
  - SemEval-2014 [Pontiki et al., 2014]
    - Restaurant and Laptop
- **Position embedding**: position embedding lookup table is initialized randomly and tuned in the training phase
- **BiGRU-PW**
  - Weights the word embeddings of each word in the sentence based on the distance from the target
- **BiGRU-PE**
  - Concatenates the word embeddings and the position embeddings of each word

| Method    | Restaurant   | Laptop       |
|-----------|--------------|--------------|
| Majority  | 65.00        | 53.45        |
| Bi-LSTM   | 78.57        | 70.53        |
| Bi-GRU    | 80.27        | 73.35        |
| Bi-GRU-PW | 79.55        | 71.94        |
| Bi-GRU-PE | 80.89        | 76.02        |
| TDLSTM    | 75.63        | 68.13        |
| MemNet    | 79.98        | 70.33        |
| IAN       | 78.60        | 72.10        |
| HAPN      | <b>82.33</b> | <b>77.27</b> |

## Aspect Sentiment Classification with both Word-level and Clause-level Attention Networks [Wang et al. 2018]

- Highlight the need for incorporating the importance of both words and clauses inside a sentence

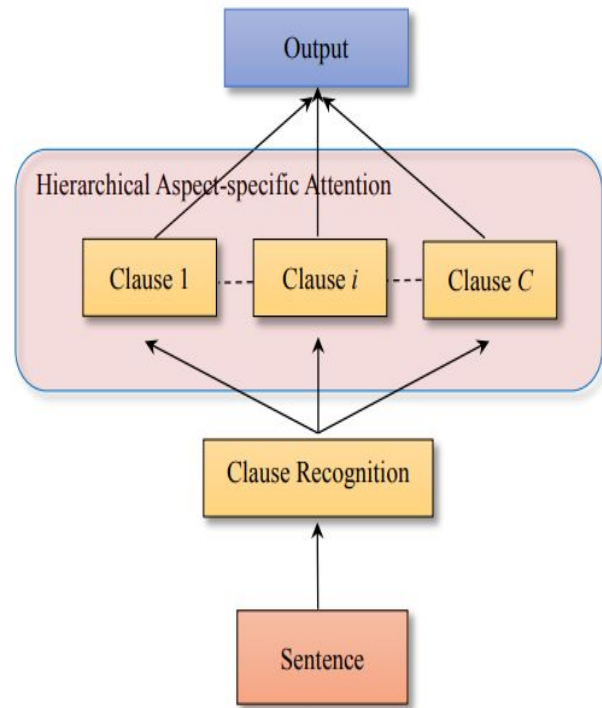


# Motivation

- For a specific aspect, importance degrees of different words are different
  - Words such as “great”, “tasty” contribute much in implying the positive sentiment polarity for the aspect FOOD#QUALITY; BUT
  - Words such as “is”, “and” don’t contribute
- For a particular aspect, the importance degrees of different clauses are different
  - the first and second clauses have much stronger information in assisting the prediction of the sentiment polarity for the aspect FOOD#QUALITY;
  - In contrast, the third and fourth clauses are more relevant to the aspect AMBIENCE#GENERAL.

# Proposed Approach

- **Clause Recognition**
  - Sentence-level discourse segmentation to segment a sentence into several clauses.
- **Hierarchical Attention**
  - **Word-level attention:** BiLSTM layers to encode all clauses and employed a word-level attention layer to capture the *importance degrees of words in each clause*
  - **Clause-level attention:** BiLSTM layer to encode the output from the former layers and propose a clause-level attention layer to capture the importance degrees of all the clauses inside a sentence



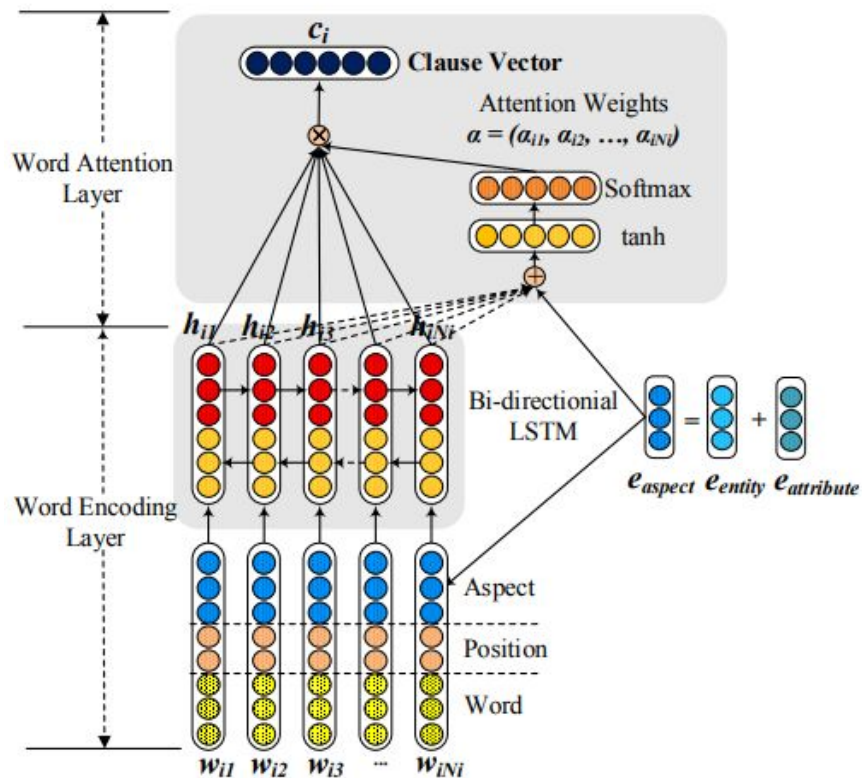


# Clause Recognition

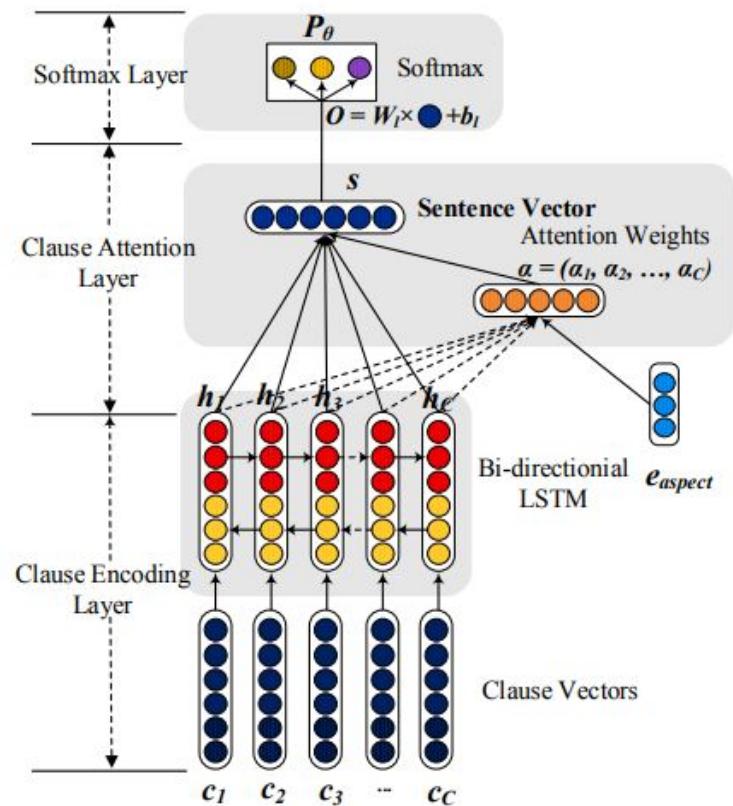
- Similar to *discourse segmentation*
  - Breaks a given text into non-overlapping segments called elementary discourse units (EDUs)
- Adopted Rhetorical Structure Theory (RST) [MANN, 1988]

[The food is great and tasty]<sup>A</sup> [and the hot dogs are especially top notch,]<sup>B</sup> [but the sitting space is too small]<sup>C</sup> [and i don't like the ambience.]<sup>D</sup>

# Architecture



(a) Word-level Aspect-specific Attention Module



(b) Clause-level Aspect-specific Attention Module

# Experiments

- Dataset

- SemEval-2015 [Pontiki et al., 2015] → Restaurant and Laptop

| Method                   | Restaurant   |              | Laptop       |              |
|--------------------------|--------------|--------------|--------------|--------------|
|                          | Accuracy     | Macro-F1     | Accuracy     | Macro-F1     |
| Majority                 | 0.537        | 0.233        | 0.570        | 0.242        |
| LSTM                     | 0.735        | 0.617        | 0.734        | 0.608        |
| TC-LSTM                  | 0.747        | 0.634        | 0.745        | 0.622        |
| ATAE-LSTM                | 0.752        | 0.641        | 0.747        | 0.637        |
| IAN                      | 0.755        | 0.639        | 0.753        | 0.625        |
| Hierarchical BiLSTM      | 0.763        | 0.647        | 0.767        | 0.632        |
| Word-level Attn          | 0.789        | 0.662        | 0.785        | 0.646        |
| Clause-level Attn        | 0.783        | 0.659        | 0.779        | 0.647        |
| Word & Clause-level Attn | <b>0.809</b> | <b>0.685</b> | <b>0.816</b> | <b>0.667</b> |

# Memory Network for Aspect Sentiment Classification

# Memory Network

- Introduced by [Weston et al. 2014]
- Core idea
  - Inference with a long-term memory component, which could be read, written to, and jointly learned with the goal of using it for prediction
- Formally,
  - A memory  $m \rightarrow$  Array of objects/vectors
  - *Four components*
    - **Input feature map ( $I$ )**  $\rightarrow$  converts input ( $x$ ) to internal feature representation
      - $I(x)$
    - **Generalization ( $G$ )**  $\rightarrow$  updates old memories with new input. Network compresses and generalizes its memories at this stage for some intended future use
      - $m_i = G(m_i, I(x), m), \forall i.$
    - **Output feature map ( $O$ )**  $\rightarrow$  generates an output representation given a new input and the current memory state,
      - $o = O(I(x), m)$
    - **Response ( $R$ )**  $\rightarrow$  outputs a response based on the output representation
      - $r = R(o).$

# Some more details: Memory Network

- **Input:** Any kind of operations possible (NER, PoS tagging, Coreference etc. on text)
- **Generalization**
  - Its main task is to store the current input in a slot of the memory
  - Update the old stored values based on the new evidence
  - Memory can be stored with topic or entity if the input is very big (Freebase, Wikipedia etc)
- **Output:** Typically responsible for reading from memory and performing inference, e.g., calculating what are the relevant memories to perform a good response
- **Response:** Produces the final response given  $O$

Example in a QA setup, *O finds relevant memories, and then R produces the actual wording of the answer, e.g., R could be an RNN that is conditioned on the output of O.*

# Why is Memory Network for ASC?

- Conventional neural models like LSTM captures context information in an *implicit way*, and are *incapable of explicitly exhibiting important context clues of an aspect*
  - Only a small subset of context words actually needed in determining the sentiment polarity
- Example: *great food but the service was dreadful!*

“dreadful” is an important clue for the aspect “service” but “great” is not needed
- Standard LSTM works in a sequential way
  - Manipulates each context word with the same operation
  - AND hence, it cannot explicitly reveal the importance of each context word

# Why is Memory Network for ASC?

- *What could be the desirable solution then?*
  - Should be capable of explicitly capturing the importance of context words
  - Use the information to build up features for the sentence after given an aspect word
- *What a human will do?*
  - will selectively focus on parts of the contexts, and
  - acquire information where it is needed to build up an internal representation towards an aspect in his/her mind

*Equivalent to store an object in memory and then search for a reasonable match*

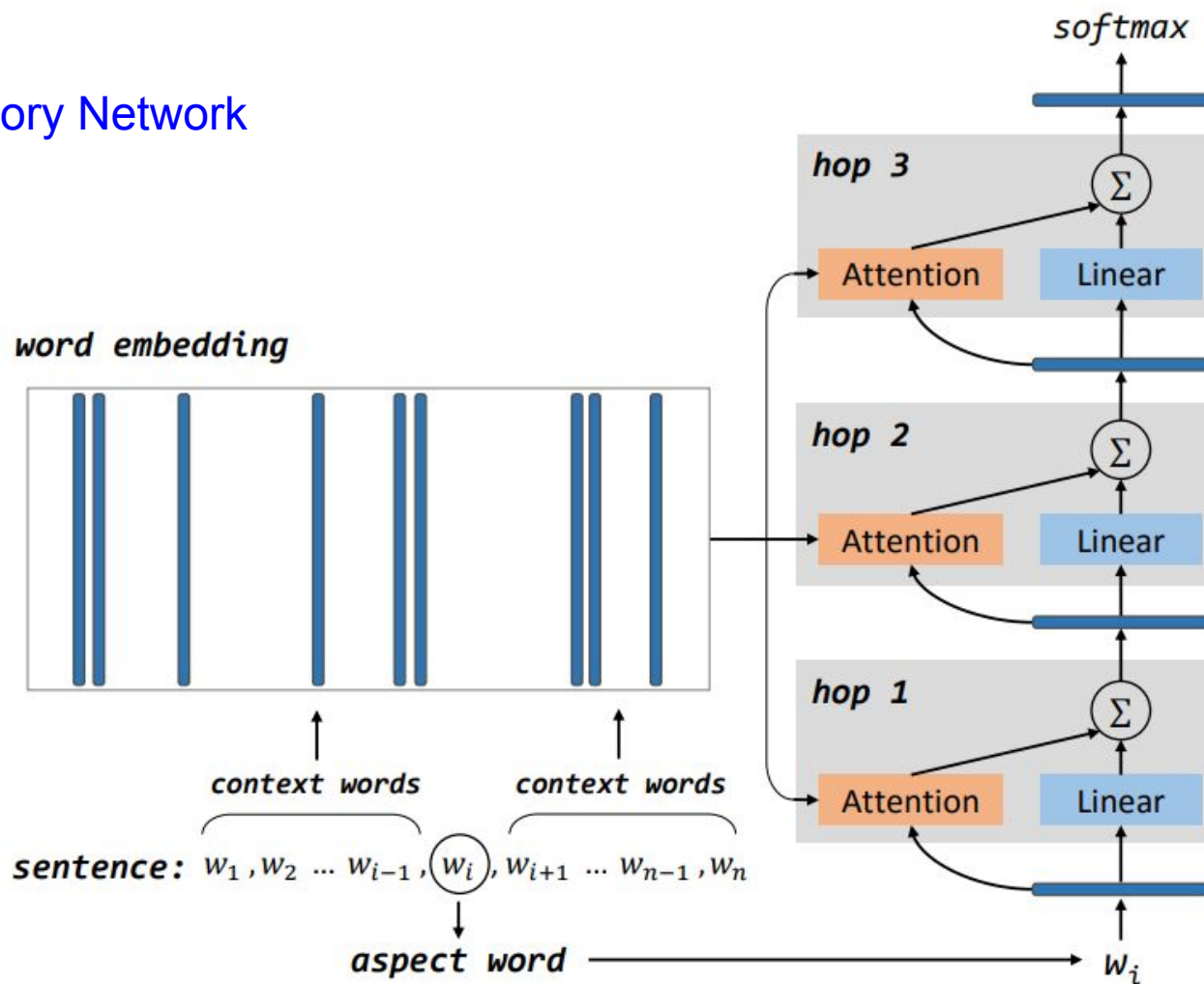


## Aspect Level Sentiment Classification with Deep Memory Network [Tang et al. 2016]

- Explicitly captures the *importance of each context word* when inferring the *sentiment polarity of an aspect*
- Utilized multiple computational layers with ***shared parameters (hops)***, each of which is a *neural attention model* over an external memory
- Each layer is a content- and location- based attention model, which *first learns the importance/weight of each context word* and then *utilizes this information to calculate the continuous text representation*

Duyu Tang, Bing Qin, Ting Liu. 2016. Aspect Level Sentiment Classification with Deep Memory Network. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 214–224, Austin, Texas, November 1-5, 2016.

## Deep Memory Network



# Attention: Content and Location

- **Content**

- Determines the most attended context word with respect to the target aspect
- Model could adaptively assign an importance score to each piece of memory according to its semantic relatedness with the aspect

- **Location**

- Sentiment-bearing word closer to the aspect is more important
- Distance of the word from the target is therefore very important
- Introduced four strategies to include the location information
- Memory content is updated based on the location attention (i.e. *how far is the memory element from the target aspect?*)

## Multiple hops

- Single attention layer is essentially a weighted average compositional function
  - Not powerful enough to handle the sophisticated computability like *negation*, *intensification* and contrary in language
- Multiple computational layers allow the deep memory network to learn representations of text with multiple levels of abstraction
- Each layer/hop retrieves important context words, and transforms the representation at previous level into a representation at a higher, slightly more abstract level

# Experiments

- Dataset
  - SemEval-2014 [Pontiki et al., 2014]

| Method              | Laptop | Restaurant |
|---------------------|--------|------------|
| Majority            | 53.45  | 65.00      |
| Feature+SVM         | 72.10  | 80.89      |
| LSTM                | 66.45  | 74.28      |
| TD-LSTM             | 68.13  | 75.63      |
| TD-LSTM + ATTENTION | 66.24  | 74.31      |
| MemNet(1)           | 67.66  | 76.10      |
| MemNet(2)           | 71.14  | 78.61      |
| MemNet(3)           | 71.74  | 79.06      |
| MemNet(4)           | 72.21  | 79.87      |
| MemNet(5)           | 71.89  | 80.14      |
| MemNet(6)           | 72.21  | 80.05      |
| MemNet(7)           | 72.37  | 80.32      |
| MemNet(8)           | 72.0   | 80.14      |
| MemNet(9)           | 72.21  | 80.95      |

## Target-Sensitive Memory Networks for Aspect Sentiment Classification [Wang et al. 2018]

- In Memory Network, attention mechanism plays a crucial role in detecting the sentiment context for the given target
- However, sentiment polarity of the (detected in memory networks) context is dependent on the given target and it cannot be inferred from the context alone
  - Sentiment contexts for both these sentences are “*high*”, i.e. the attention mechanism will have higher weights for “*high*” in both the cases
    - *The **price** is high.* → Negative
    - *The **screen resolution** is high.* → Positive
- Incorporate target information to infer **(price, high) as negative** and **(screen resolution, high) as positive**

## Six variants of TMNs

### 1. Non-linear Projection (NP)

- $\alpha_i$ ,  $c_i$  and  $v_t \Rightarrow$  Attention score, Context and vector of target ( $t$ )
- Interaction between target and context

$$s = W \cdot \tanh\left(\sum_i \alpha_i c_i + v_t\right)$$

### 2. Contextual Non-linear Projection (CNP)

$$s = W \sum_i \alpha_i \cdot \tanh(c_i + v_t)$$

### 3. Interaction Term (IT):

- It measures the sentiment-oriented interaction effect between targets and contexts, i.e., Target-Context-Sentiment (TCS)

$$s = \sum_i \alpha_i (W_s c_i + w_I \langle d_i, d_t \rangle)$$
$$d_i = D x_i, d_t = D t$$

$D$  = Embedding matrix that captures the sentiment interactions

## Six variants of TMNs

### 4. **Coupled Interaction (CI):**

- Additionally captures the global correlation between context and different sentiment classes.

$$s = \sum_i \alpha_i (W_s c_i + W_I \langle d_i, d_t \rangle e_i)$$

### 5. **Joint Coupled Interaction (JCI):**

- Simplification of CI model

$$s = \sum_i \alpha_i (W_s c_i + W_I \langle d_i, d_t \rangle c_i)$$

### 6. **Joint Projected Interaction (JPI)**

- First component, captures target-independent sentiment effect
- Second component, TCS interaction

$$s = \sum_i \alpha_i W_J \tanh(W_1 c_i) + \sum_i \alpha_i W_J \langle d_i, d_t \rangle \tanh(W_2 c_i)$$



# Experiments

- Dataset: SemEval-2014 [Pontiki et al., 2014]
  - Restaurant and Laptop

| Method           | Restaurant |       | Laptop |       |
|------------------|------------|-------|--------|-------|
|                  | 1-hop      | 3-hop | 1-hop  | 3-hop |
| <b>AE-LSTM</b>   | 66.45      | -     | 62.45  | -     |
| <b>ATAE-LSTM</b> | 65.41      | -     | 59.41  | -     |
| <b>NP</b>        | 64.62      | 65.98 | 62.63  | 67.79 |
| <b>CNP</b>       | 65.58      | 66.87 | 64.38  | 64.85 |
| <b>IT</b>        | 65.37      | 68.64 | 63.07  | 66.23 |
| <b>CI</b>        | 66.78      | 68.49 | 63.65  | 66.79 |
| <b>JCI</b>       | 66.21      | 68.84 | 64.19  | 67.23 |
| <b>JPI</b>       | 66.58      | 67.86 | 64.53  | 64.16 |

## IARM: Inter-Aspect Relation Modeling with Memory Networks in Aspect-Based Sentiment Analysis [Majumder et al. 2018]

- Incorporates the neighboring aspects related information for sentiment classification of the target aspect (i.e. *there is a dependency between different aspect terms*)
  - Example 1: “*The **menu** is very limited - I think we counted 4 or 5 **entries**.*”
    - Non-trivial to predict the sentiment for aspect “**entries**”, unless the other aspect “**menu**” is considered
    - Negative sentiment of “**menu**” induces “**entries**” to have the same sentiment
  - Example 2: “***Food** is usually very good, though I wonder about freshness of **raw vegetables***”
    - No clear sentiment marker for “**raw vegetables**”
    - The **positive** sentiment of “**food**”, due to the word “**good**”, and the presence of conjunction “**though**” determines the sentiment of “**raw vegetables**” to be **negative**

# Method: Key Steps

- **Input Representation**

- Input sentences and aspect-terms are represented using pre-trained Glove word embeddings
- For multi-worded aspect-terms, we take the mean of constituent word embeddings as aspect representation

- **Aspect-Aware Sentence Representation**

- Embedding of each word in a sentence is concatenated with the given aspect representation
- Modified sequence of words is fed to a GRU for context propagation
- Attention layer to obtain the aspect-aware sentence representation (for all the aspects in a sentence)

# Method: Key Steps

- **Inter-Aspect Dependency Modeling**
  - Match the target-aspect-aware sentence representation with aspect-aware sentence representation of the other aspects
  - More refined sentence representation after a certain number of iterations of the memory network
  - Softmax layer for final classification

## Inter-Aspect Relation Modeling (IARM)

- Aspect-Aware Sentence Representation (AASR)  
[Wang et al. 2016]

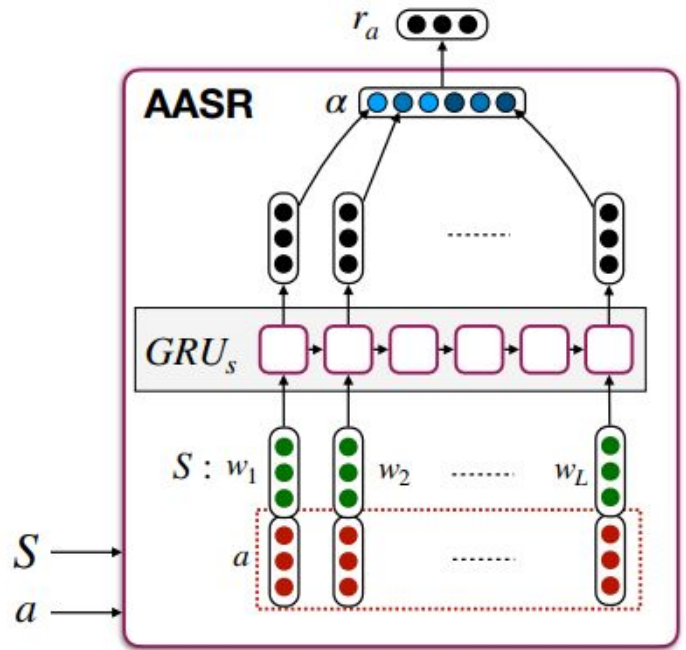
$$S_{ai} = [w_1 \oplus a_i, w_2 \oplus a_i, \dots, w_L \oplus a_i] \in \mathbb{R}^{L \times 2D}$$

$$R_{ai} = GRU(S_{ai})$$

$$\alpha = \text{softmax}(R_{ai} W_s + b)$$

$$r_{ai} = \alpha^T R_{ai}$$

$$R = [r_{a_1}, r_{a_2}, \dots, r_{a_M}]$$



Aspect-aware Sentence Representation

# Inter-Aspect Relation Modeling (IARM)

- Inter-Aspect Dependency Modeling
  - Models the dependency of the target aspect with the other aspects in the sentence.

$$Q = GRU_a(R)$$

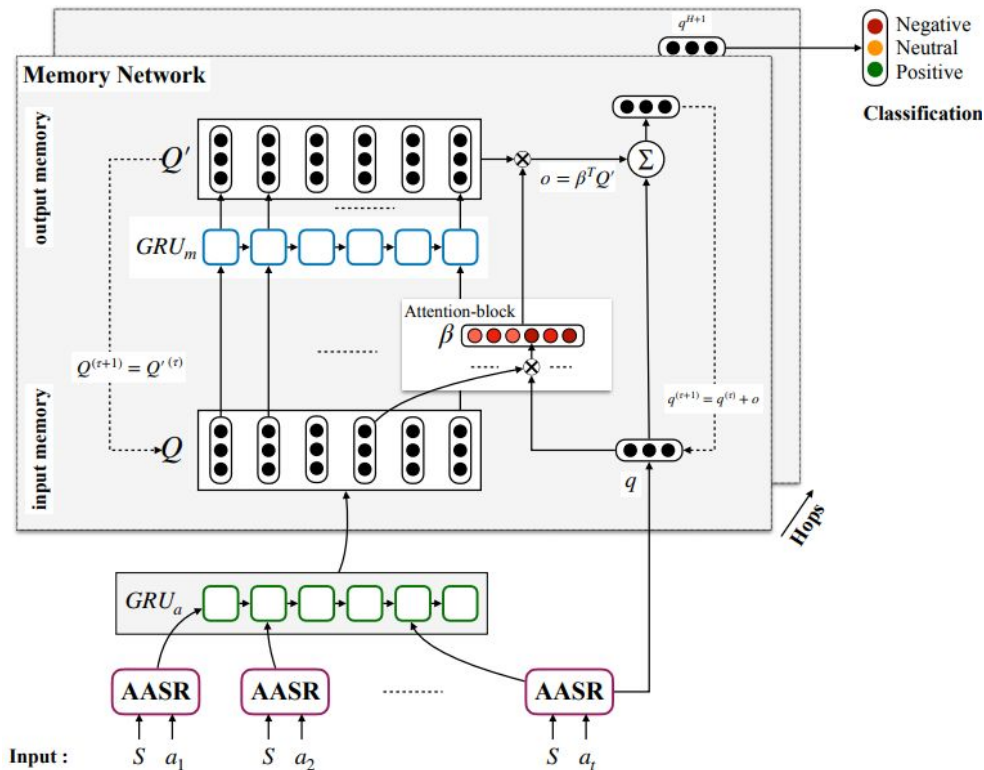
$$q = \tanh(r_{a_t} W_T + b_T)$$

$$z = qQ^T,$$

$$\beta = \text{softmax}(z)$$

$$Q' = GRU_m(Q)$$

$$o = \beta^T Q'$$



## Experimental Results

- Dataset: SemEval-2014 [Pontiki et al., 2014]
  - Restaurant and Laptop

| Method    | Restaurant   | Laptop       |
|-----------|--------------|--------------|
| Majority  | 0.535        | 0.650        |
| LSTM      | 0.743        | 0.665        |
| TD-LSTM   | 0.756        | 0.681        |
| AE-LSTM   | 0.762        | 0.689        |
| ATAE-LSTM | 0.772        | 0.687        |
| IAN       | 0.786        | 0.721        |
| IARM      | <b>0.800</b> | <b>0.738</b> |

# Attention weights: Heatmaps

Example 1:

*"I recommend any of their **salmon dishes**."*



(a) Attention weight for aspect "salmon dishes" for IAN.



(b) Attention weight for aspect "salmon dishes" for IARM.



(a) Attention weights for aspect "cosi sandwiches" for IAN.



(b) Attention weights for aspect "cosi sandwiches" for IARM.



(c) Attention weights for aspect "coffee" for IARM.

Example 2:

*"**Coffee** is a better deal than overpriced **cosi sandwiches**."*



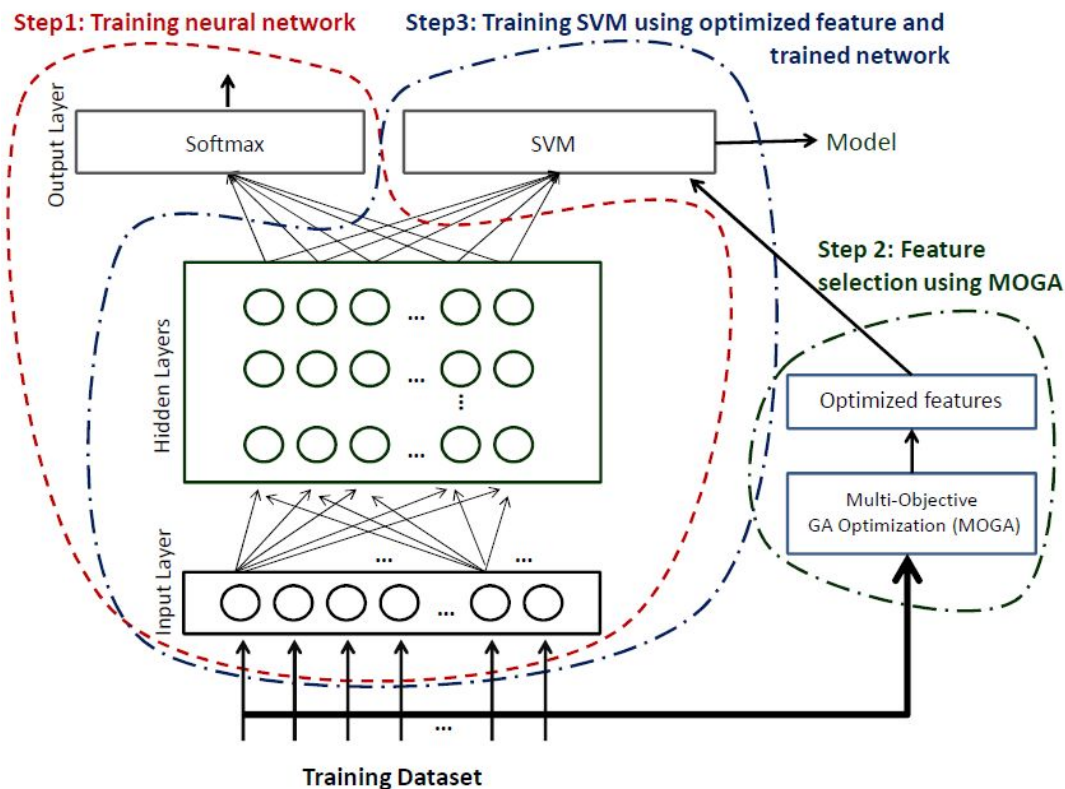
***CNN + Hand-crafted Optimized Features + SVM***

# A Hybrid Deep Learning Architecture for Sentiment Analysis [Akhtar et al. 2016]

- CNN based hybrid architecture for sentiment analysis
  - Replace a weak classifier (**softmax regression**) with a stronger classifier (**SVM**) at the output layer
- Assist CNN with optimized feature set obtained through GA based multiobjective optimization
- For each aspect, look for the sentiment marker near the aspect term itself
  - Define context as +/- few words (e.g., 3) in the neighbourhood, i.e., 3 prev tokens and 3 next tokens
    - ***Tech support*** would not fix the problem
      - [null, null, null, ***Tech\_support***, would, not, fix]
    - ***The entire place*** is very clean
      - [null, the, entire, ***place***, is, very, clean]

# Classification Model

1. Training of a typical convolutional neural network (CNN)
  - Obtain weight matrix
2. A multi-objective GA based optimization technique (NSGA-II) for extracting the optimized set of features
  - Two objectives
    - *Accuracy* (maximize)
    - *Num of features* (minimize)
3. Training of SVM utilizing the network trained in first step and optimized features



# Datasets

- **Hindi**

- Twitter (SAIL - 2015): 1.6K sentences
- Product/Service reviews: 5.4K sentences
  - Aspect based sentiment analysis
  - Sentence based sentiment analysis
- Movie reviews: 2.1K sentences
  - Sentence based sentiment analysis

- **English**

- Twitter (SemEval - 2015): 10.2K sentences
  - Generic tweets
  - Sarcastic tweets
- Product/Service reviews (SemEval - 2014): 7.6K sentences
  - Aspect based sentiment analysis

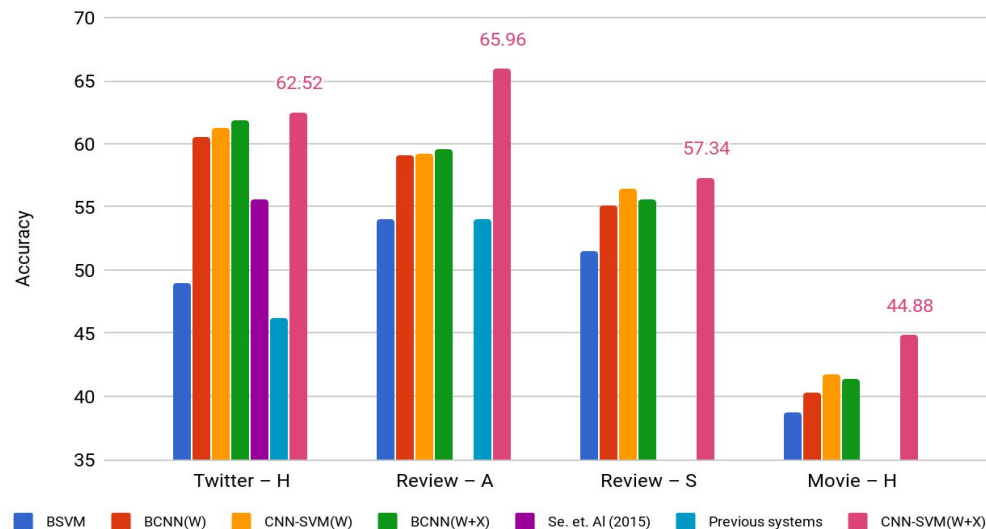
Hindi product and movie reviews datasets are available at: <http://www.iitp.ac.in/~ai-nlp-ml/resources.html>

## Feature set

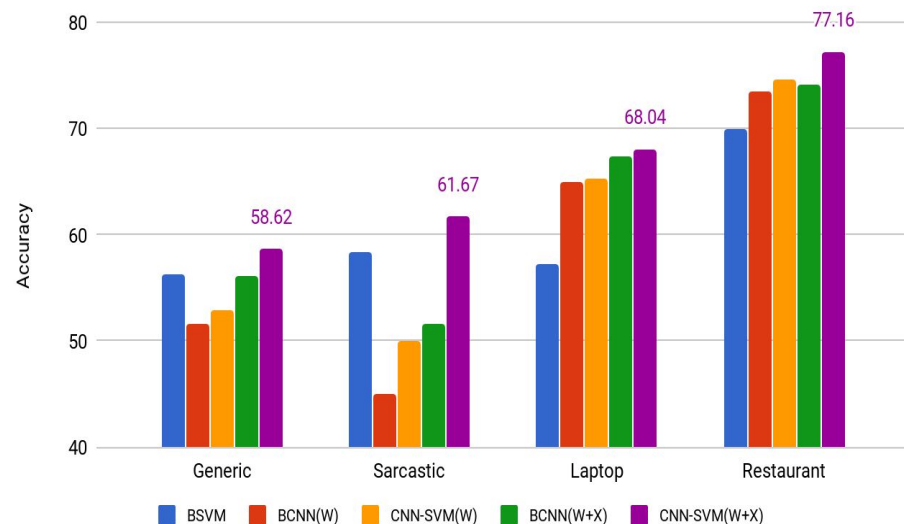
| Language | Dataset                | Optimized Features (NSGA-II)                  |
|----------|------------------------|---|
| Hindi    | Twitter                | Emoticons, Punctuation, SentiWordNet          |
|          | Review – A, Review – S | Semantic Orientation (SO)                     |
|          | Movie                  | Semantic Orientation (SO), SentiWordNet       |
| English  | Twitter                | Hashtag, Emoticons, Punctuation, BingLiu, NRC |
|          | Review – A             | BingLiu, MPQA                                 |

# Evaluation

## Hindi Datasets



## English Datasets



- $B_{SVM}$  : SVM based model
- $B_{CNN(W)}$  : CNN based model with word vectors as input
- $B_{CNN(W+X)}$  : CNN based model with word vectors and optimized feature set as input.
- $CNN-SVM^{(W)}$  : SVM on top of CNN with word vectors as input.
- $CNN-SVM^{(W+X)}$  : SVM on top of CNN with word vectors and optimized feature set as input.

## *Cross-lingual and Multi-lingual ASC*

## Solving Data Sparsity for Aspect based Sentiment Analysis using Cross-linguality and Multi-linguality [Akhtar et al. 2018]

- Low-resource languages usually suffer in performance due to the *non-availability of sufficient* training data instances.
- Low-resource languages (e.g. Hindi, Bengali etc.) usually suffer due to the non-availability of sufficient data instances
- **Problem:** Data Sparsity in word representation (i.e. absence of representation of a word) is another problem
- Out-of-vocabulary (OOV) words in a word embedding model pose a serious challenge to the underlying learning algorithm

Shad Akhtar, Palaash Sawant, Sukanta Sen, Asif Ekbal, and Pushpak Bhattacharyya (2018). Solving Data Sparsity for Aspect based Sentiment Analysis using Cross-linguality and Multi-linguality. In Proceedings of the 16th Annual Conference of the NAACL:HLT-2018, June 2018, New Orleans, LAUSA, pages 572–582.



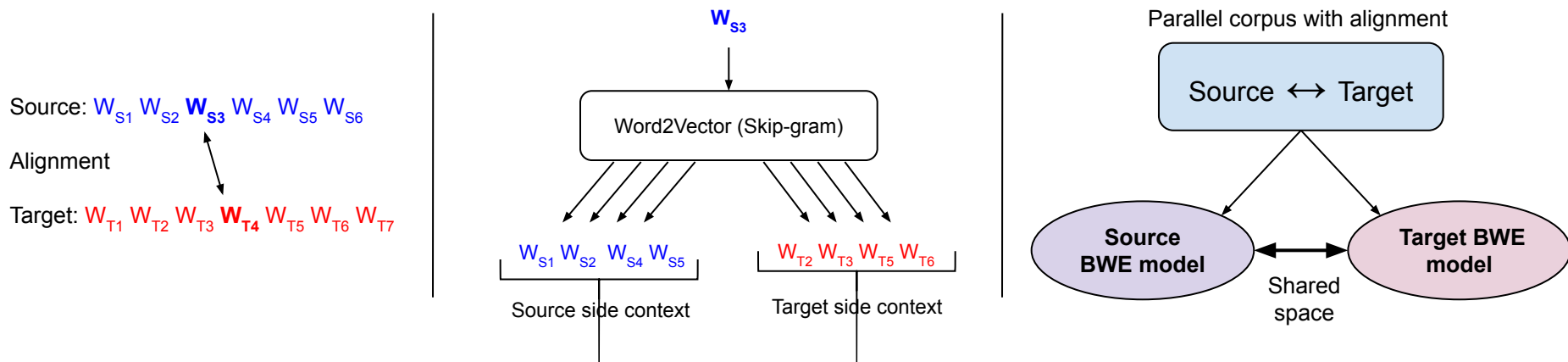
# Solution to the OOV problem

- **Solution:** Minimize the effect of data sparsity problem in a resource-scarce language scenario by leveraging the information of resource-rich languages
  - How?
    - Word embedding space of two languages may not be same
    - Therefore, cannot use the two embeddings in the similar context
  - Project the embeddings of two languages into a shared space
    - Bi-lingual embeddings (Luong et al., 2015)

Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. 2015, ***Bilingual Word Representations with Monolingual Quality in Mind***. In *NAACL Workshop on Vector Space Modeling for NLP*.

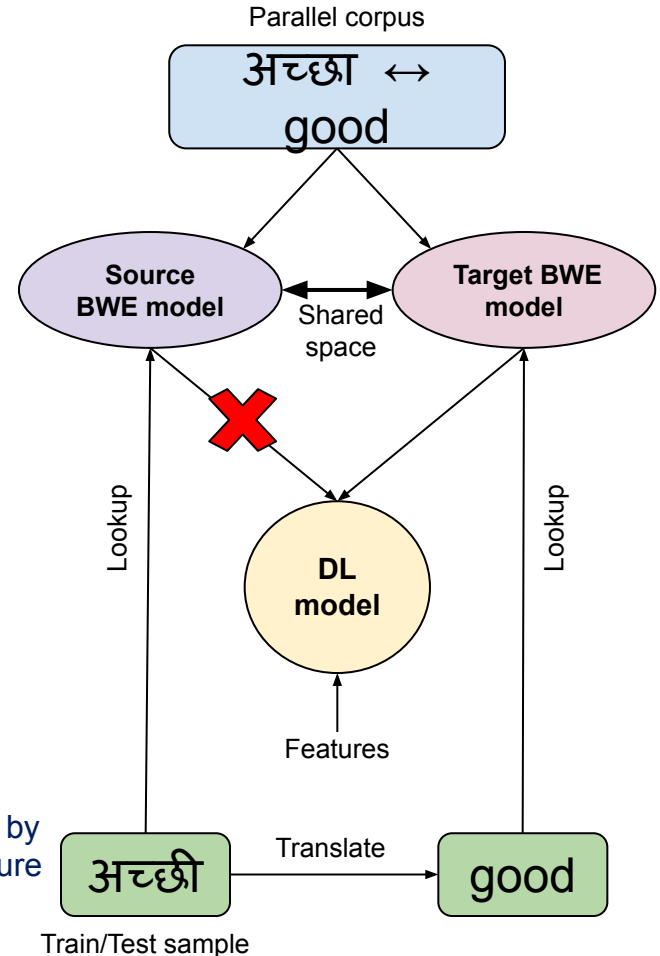
# Bi-lingual Word Embeddings (BWE): (Loung et al., 2015)

- Bi-lingual word embeddings aims to *bridge the language divergence in the vector space*
  - Requires a *parallel corpus* and *alignment information* among parallel sentences
  - Utilize existing word2vec skip-gram model (Mikolov et al., 2013)
  - For each word, the authors defined its context to include the neighbouring words from both the source and target languages



# Proposed Approach

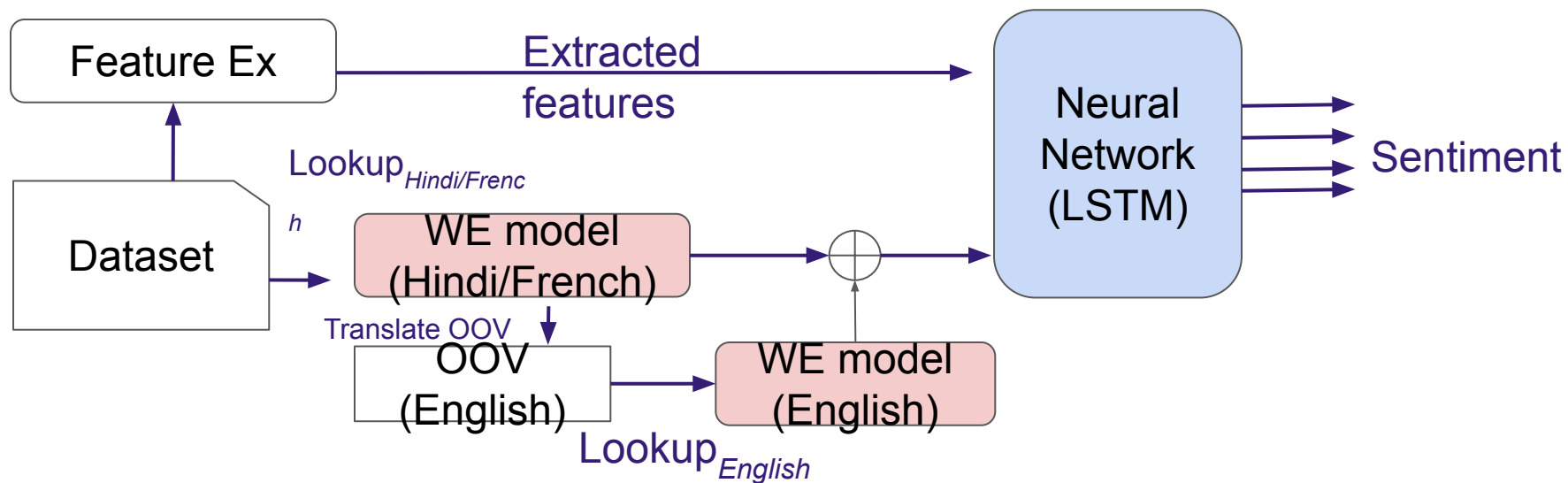
- Utilize bilingual word embeddings for a pair of languages (e.g., Hindi-English, French-English)
- Extract word representations for all the words in a sentence from the source (e.g., Hindi) bilingual word embedding
- For all the OOV words, translate into the target language, and perform another lookup in target bilingual embedding
- Spelling variation: Two differently spelled words in Hindi such as 'किबनशन | *kambineshana*' and 'कंबीनशन | *kaMbIneshana*' translate to an English word "combination"
- Further, leverage the effectiveness of English side resources by translating a word into English and then extracting its feature representation
  - Bing Liu, MPQA, SentiWordNet and Semantic Orientation



# Two setups

- Multi-lingual Setup
  - **Train** and **Test** on **Source** language (i.e., Hindi or French)
  - Utilize bi-lingual embeddings for OOV words
  - Utilize English-side lexicons for the feature extraction
- Cross-lingual Setup
  - **Train** on **Target** language (i.e., English) and **Test** on **Source** language (i.e., Hindi or French)
  - Utilize bi-lingual embeddings for OOV words
  - Utilize English-side lexicons for the feature extraction

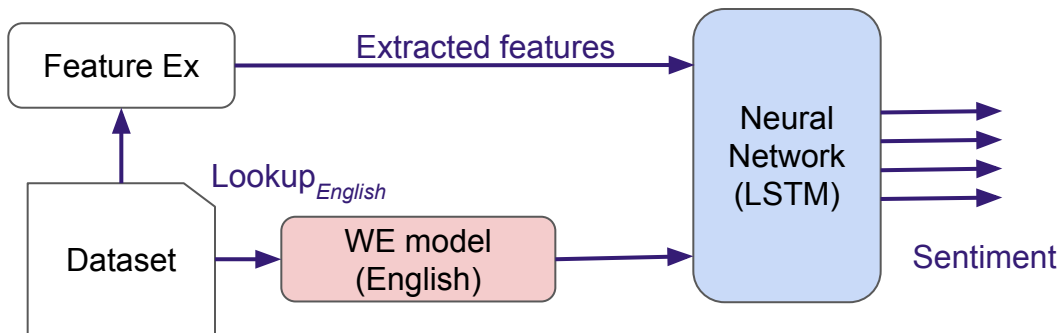
# Multi-lingual Setup



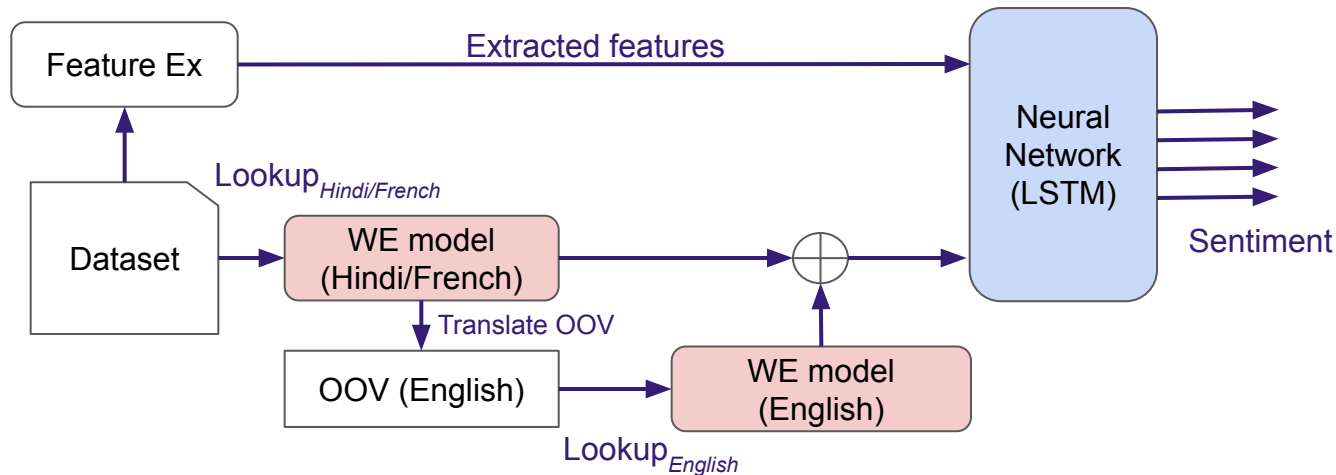
**Training and Testing scenarios**

# Cross-lingual Setup

Training



Testing



# Hybrid Architecture

Three architectures based on the position of the fusion of hand-crafted features

## A1. Early fusion:

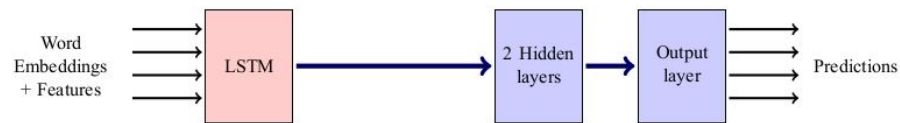
- $\text{LSTM}(WE + Feat)$

## A2. Delayed fusion:

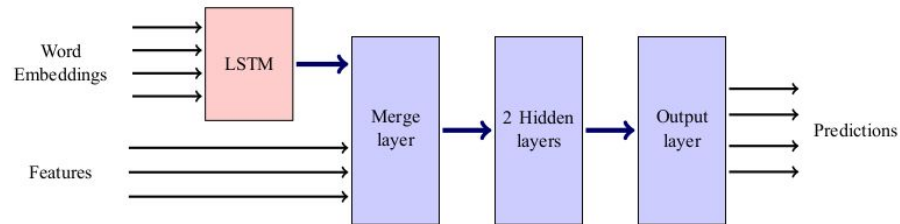
- $\text{LSTM}(WE) + Feat$

## A3. Delayed fusion with sequential feature representation

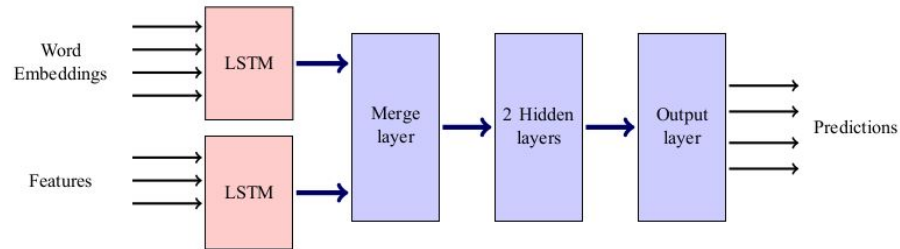
- $\text{LSTM}(WE) + \text{LSTM}(Feat)$



(a) Architecture A1



(b) Architecture A2



(c) Architecture A3

# Dataset and Experimental setups

- Aspect Based Sentiment Analysis

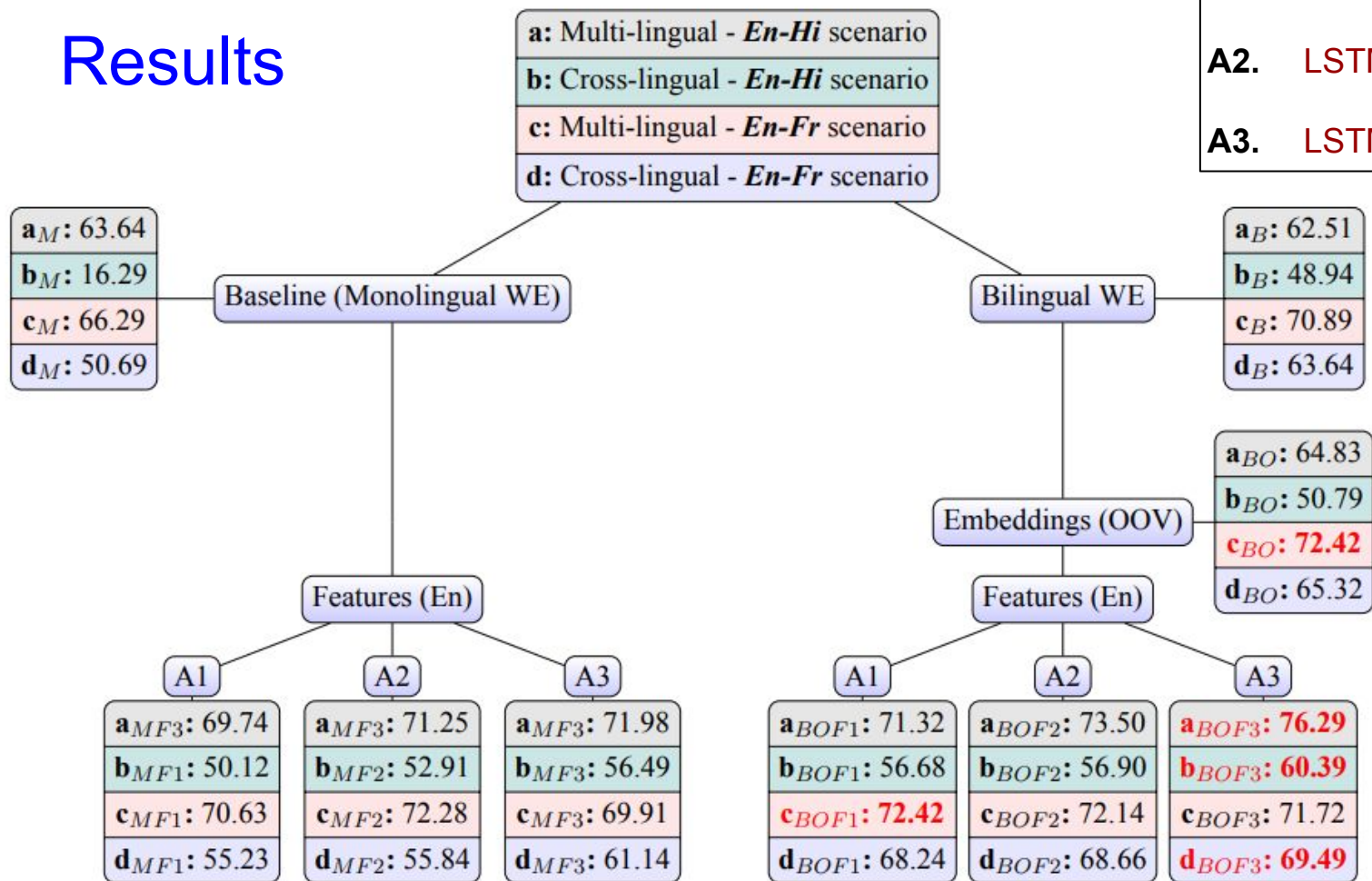
| Language pairs   | Datasets                                      | Review Sentences | Aspect terms |
|------------------|---|------------------|--------------|
| English - Hindi  | English - SemEval 2014 (Pontiki et al., 2014) | 3845             | 3012         |
|                  | Hindi (Akhtar et al., 2016)                   | 5417             | 4509         |
| English - French | English - SemEval 2016 (Pontiki et al., 2016) | 3365             | 2676         |
|                  | French - SemEval 2016 (Pontiki et al., 2016)  | 2429             | 3482         |

- Setups

- Multi-lingual Setup
  - *Train and Test on Source language (i.e., Hindi or French)*
  - Utilize English-side lexicons for the feature extraction.
- Cross-lingual Setup
  - *Train on Target language (i.e., English) and Test on Source language (i.e., Hindi or French)*
  - Utilize English-side lexicons for the feature extraction.



# Results



\*Accuracy values

# Summary and Takeaways

- **Summary**

- Presented the background of ABSA
- Presented the state-of-the-art deep learning models like LSTM, LSTM with attention, GRU, Memory networks etc. for aspect classification

- **Takeaways**

- LSTM with target-specific attention helps obtaining good accuracy for ASC
- Encoding position of the aspect term in the sentence helps for better classification
- Interactive attention (aspect-aware as well as context-aware representations) can better disambiguate the classification
- Hierarchical attention (attention at aspect level to find most matching aspect term + attention to find the best sentiment bearing words) is useful
- Memory network could be employed to model the inter-aspect relations
- Cross-lingual embedding representation is important to perform multi-lingual and cross-lingual SA involving low-resource languages

# Future Works

- Sentiment intensity prediction in ABSA
- ABSA in multi-modal scenario
- Effective solutions to ABSA in low-resource scenario
  - Cross-lingual embedding
  - Injecting external knowledge base into deep neural network
  - Transfer learning and domain adaptation

## References

- Zeiler, Matthew D. and Fergus, Rob (2014). Visualizing and Understanding Convolutional Networks. Computer Vision -- ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I. Springer International Publishing
- Sepp Hochreiter and Jurgen Schmidhuber. 1997. "Long short-term memory. Neural computation, 9(8):1735–1780.
- Duyu Tang, Bing Qin, Xiaocheng Feng, Ting Liu. 2016. Effective LSTMs for Target-Dependent Sentiment Classification. In Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, pages 3298–3307, Osaka, Japan, December 11-17 2016.
- Md Shad Akhtar, Palaash Sawant, Sukanta Sen, Asif Ekbal, and Pushpak Bhattacharyya (2018). Solving Data Sparsity for Aspect based Sentiment Analysis using Cross-linguality and Multi-linguality. In Proceedings of the 16th Annual Conference of the NAACL:HLT-2018, June 2018, New Orleans, LA, USA, pages 572–582.
- Yequan Wang, Minlie Huang, Li Zhao and Xiaoyan Zhu. 2016. Attention-based LSTM for Aspect-level Sentiment Classification. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 606–615, Austin, Texas, November 1-5, 2016.
- Jiajun Cheng, Shenglin Zhao, Jiani Zhang, Irwin King, Xin Zhang, Hui Wang. 2017. Aspect-level Sentiment Classification with HEAT (HiErarchical ATtention) Network. In Proceedings of the CIKM-17, November 6-10, 2017, Singapore.
- Dehong Ma, Sujian Li, Xiaodong Zhang, Houfeng Wang. 2017. Interactive Attention Networks for Aspect-Level Sentiment Classification. Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI-17).
- Ruidan He, Wee Sun Lee, Hwee Tou Ng, and Daniel Dahlmeier. 2018. Effective Attention Modeling for Aspect-Level Sentiment Classification. In Proceedings of the 27th International Conference on Computational Linguistics, pages 1121–1131 Santa Fe, New Mexico, USA, August 20-26.
- Lishuang Li, Yang Liu and AnQiao Zhou. 2018. Hierarchical Attention Based Position-aware Network for Aspect-level Sentiment Analysis. In Proceedings of the 22nd Conference on Computational Natural Language Learning (CoNLL2018), pages 181-189 Brussels, Belgium, Oct 31-Nov 1.
- Jingjing Wang, Jie Li, Shoushan Li, Yangyang Kang, Min Zhang, Luo Si, Guodong Zhou. 2018. Aspect Sentiment Classification with both Word-level and Clause-level Attention Networks. In Proceedings of the 22nd Conference on Computational Natural Language Learning (CoNLL2018), pages 181-189 Brussels, Belgium, Oct 31-Nov 1.
- XINYI WANG, GUANGLUAN XU, JINGYUAN ZHANG, XIAN SUN, LEI WANG, AND TINGLEI HUANG. 2018. Syntax-Directed Hybrid Attention Network for Aspect-Level Sentiment Analysis. IEEE Access, Jan. 2019.
- Duyu Tang, Bing Qin, Ting Liu. 2016. Attention-based LSTM for Aspect-level Sentiment Classification. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 214–224, Austin, Texas, November 1-5, 2016.
- Cheng Li, Xiaoxiao Guo, Qiaozhu Mei. 2017. Deep Memory Networks for Attitude Identification. In Proceedings of the WSDM 2017, February 06-10, 2017, Cambridge, United Kingdom.
- Shuai Wang, Sahisnu Mazumder, Bing Liu, Mianwei Zhou, Yi Chang. 2018. Target-Sensitive Memory Networks for Aspect Sentiment Classification. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Long Papers), pages 957–967 Melbourne, Australia, July 15 - 20, 2018.
- Navonil Majumder, Soujanya Poria, Alexander Gelbukh, Md Shad Akhtar, Erik Cambria, Asif Ekbal. 2018. IARM: Inter-Aspect Relation Modeling with Memory Networks in Aspect-Based Sentiment Analysis. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 3402–3411 Brussels, Belgium, October 31 - November 4, 2018.

***Thank you for your  
attention!***

*A Few more stuffs*

## Syntax-Directed Hybrid Attention Network for Aspect-Level Sentiment Analysis [Wang et al. 2019]

- (Global) Attention mechanism that attends to all words in the context to model the interaction between target and sentence *suffers* from assigning high-attention score to irrelevant sentiment words
  - “The **wait staff** is very friendly, if you are not rude or picky”.
  - ‘Rude’ may get unnecessary high score for the target ‘wait staff’.
- Further, position vector may not work if the sentiment-oriented word is very far from the target
  - “Apple is unmatched in product quality, aesthetics, craftsmanship, and **customer service**”
  - The target ‘customer service’ is distant apart from the word ‘unmatched’.

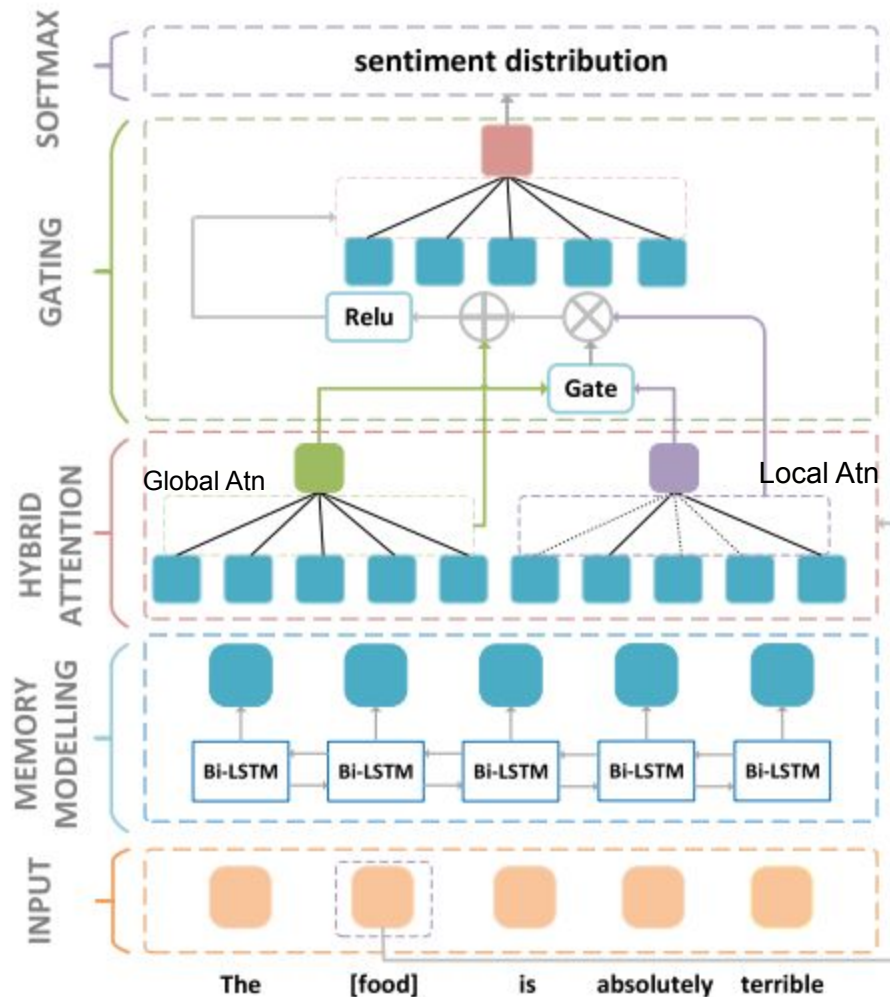
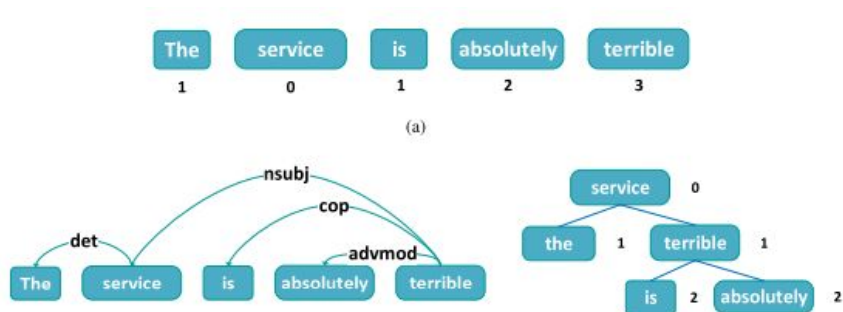
## Syntax-directed hybrid attention network (SHAN)

- A *global attention* is employed to capture coarse information about the target
- A *syntax-directed local attention* is used to take a look at words syntactically close to the target
- Utilizing global and local attention information, a less-noisy and more sentiment-oriented representation is obtained



# Architecture

- Global Attention over entire sentence
- Local Attention over syntactically closer words only
  - For target 'service', 'terrible' is closer than 'absolutely'



# Experiments

- Dataset
  - SemEval-2014 [Pontiki et al., 2014]
    - Restaurant and Laptop

| Method             | Restaurant   | Laptop       |
|--------------------|--------------|--------------|
| <b>SVM</b>         | 80.16        | 70.49        |
| <b>LSTM</b>        | 74.30        | 66.50        |
| <b>TDLSTM</b>      | 75.63        | 68.13        |
| <b>AT-LSTM</b>     | 77.20        | 68.90        |
| <b>IAN</b>         | 78.60        | 72.10        |
| <b>BiLSTM-Attn</b> | 78.16        | 72.41        |
| <b>HEAT-BiGRU</b>  | 78.68        | 73.17        |
| <b>MemNet</b>      | 78.16        | 70.33        |
| <b>SHAN</b>        | <b>81.02</b> | <b>74.64</b> |

## Deep Memory Networks for Attitude Identification [Li et al. 2017]

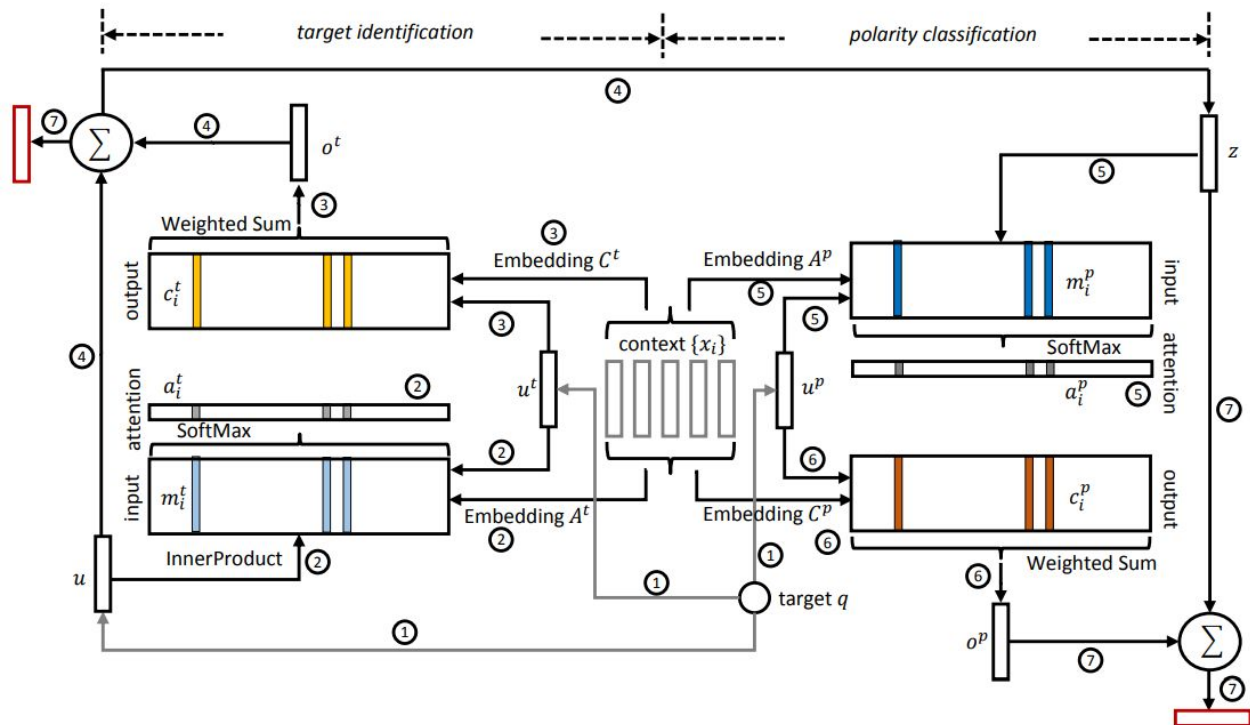
- Sentiment vs. Opinion vs. Attitude
  - I feel happy → Sentiment without target
  - We should do more exercise → Opinion without polarity
  - IJCAI is a great conference → Attitude (Polarity) towards a particular entity
- Two subtasks
  - Target Detection
    - Target can be implicit and explicit
      - “*We have been waiting for food for one hour.*” → Target is service not food
  - Polarity classification

## Deep Memory Networks for Attitude Identification [Li et al. 2017]

- Proposed a **joint framework for target word detection (TD) and polarity classification (PC)**
- Signals identified in the first subtask – both the words that refer to the target and the positions of these words, could provide useful information for the polarity of sentiments
  - “*this **camera** is \_\_\_\_\_*” → Information about the target indicates that the blank space could be **flavour or price**
  - “*this \_\_\_\_\_ is **awesome***” → The sentiment expression signal the existence of a **target**

# Attitude Network (AttNet)

1. Target Embedding
2. Input Representation and Attention for TD
3. Interleaving TD and PC
4. Input Representation and Attention for PC
5. Output Representation for PC
6. Prediction for TD and PC



# Experiments

- Dataset: SemEval-2014 + SemEval-2015 [Pontiki et al., 2014, 2015]
  - Restaurant and Laptop

sep: Separate model for both TD and PC  
sgl: Single model for TD and PC

| Method                 | F-score | Precision | Recall |
|------------------------|---------|-----------|--------|
| <b>SVM-sep</b>         | 38.43   | 51.22     | 36.83  |
| <b>SVM-sgl</b>         | 36.06   | 50.79     | 34.07  |
| <b>CNN-sep</b>         | 37.15   | 43.73     | 33.24  |
| <b>CNN-sgl</b>         | 35.45   | 44.65     | 32.83  |
| <b>BiLSTM-sep</b>      | 40.78   | 42.54     | 39.01  |
| <b>BiLSTM-sgl</b>      | 39.68   | 41.88     | 38.81  |
| <b>MultiBiLSTM-sep</b> | 40.47   | 44.89     | 37.67  |
| <b>MultiBiLSTM-sgl</b> | 39.38   | 43.22     | 37.92  |
| <b>MemNet-sep</b>      | 41.75   | 45.61     | 39.25  |
| <b>MemNet-sgl</b>      | 41.65   | 45.23     | 39.13  |
| <b>AttNet</b>          | 45.93   | 50.34     | 44.95  |