



# Arabic Social Media Sentiment Analysis Using RCNN

## Workshop

**Rania Abd EL Monam Mohamed**

*Department of Computer Sciences  
Faculty of Graduate Studies for Statistical Research  
Cairo University, Egypt*

11/12/2019

# Agenda

1. Introduction.
2. Proposed Corpus.
3. Embedding Layer.
4. Architecture of RCNN Model.
5. Publications.
6. RCNN Code.

# Agenda

1. Introduction.
2. Proposed Corpus.
3. Embedding Layer.
4. Architecture of RCNN Model.
5. Publications.
6. RCNN Code.

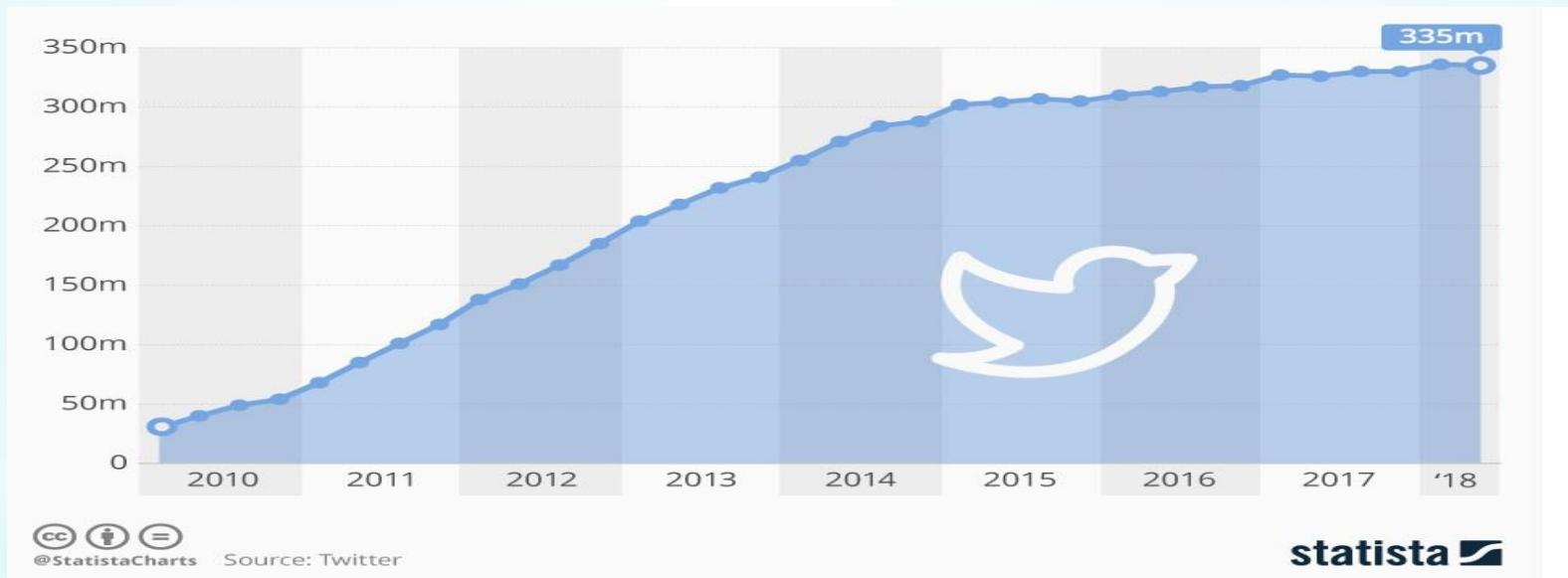
# Social Media

- Social Media is generating a large volume of sentiment rich data in the form of tweets, status updates, blog posts, comments, reviews, etc.
- Social Media are considered an excellent source of information and can provide opinions, thoughts, and insights toward various important topics.
- It is expected that the estimated amount of data on the web will be about 40 trillion gigabytes, in 2020 [1].

---

1- Gantz, J. and D. Reinsel, *Digital Universe Study*. Extracting Value from Chaos, IDC Go-to-Market Services, 2011.

- ❑ Twitter is one of the most popular, widely used social media.
- ❑ Official Twitter statistics show that 335 million monthly users actively used Twitter in 2018.



Number of monthly active twitter users around the world from 2010 to 2018 [2]

---

2- Twitter, <https://www.businessofapps.com/data/twitter-statistics/>. 2018.

# Sentiment Analysis

- ❑ Sentiment Analysis (SA) also known as **review mining**.
- ❑ SA can be defined as a process that **automates mining** of attitudes, opinions and emotions from text.
- ❑ SA is the process of determining whether **a piece of writing** is positive 😊 or negative 😞.
- ❑ SA refers to the process of deriving **high-quality information** from text.

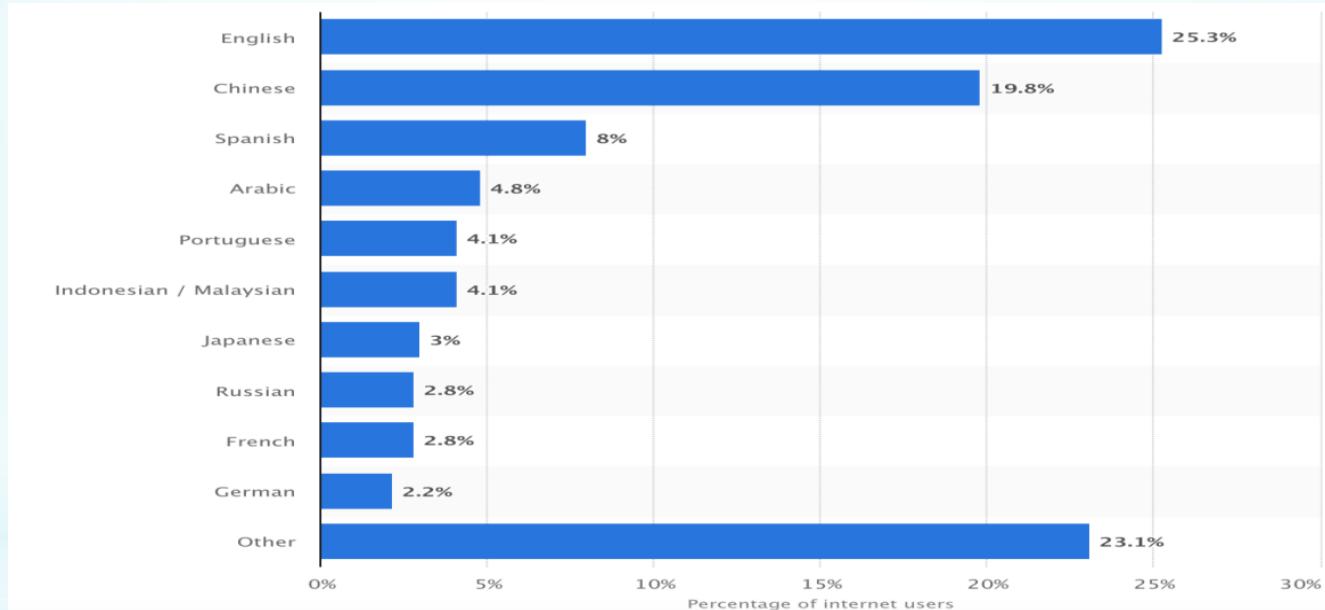
- ❑ SA plays a significant role in our **daily decision making** process.
- ❑ SA **helps** in achieving various **goals** like observing public mood regarding political topics, market intelligence, new product sales prediction, the measurement of customer satisfaction and many more.
- ❑ SA can be divided into **three levels**, namely
  - ✓ Document Level
  - ✓ Sentence Level.
  - ✓ Aspect Level (also known as word- or feature- level).

# Arabic Contents

- ❑ Arabic language is one of **six official** languages of the United Nations.
- ❑ Arabic is the official language of **21 countries**, and it's the major language in several areas of the world.
- ❑ Arabic language is classified into **three** types;
  - ✓ Classical Arabic ("AL-Qur'an")
  - ✓ Modern Standard Arabic (formal communications, television, radio, news, education)
  - ✓ Colloquial Arabic (shopping, chatting or in their homes.)

# Arabic Contents

- Arabic is currently ranked as the **forth** language used in the web, and there are about 168 million of Arabic Internet user [5].



Top 10 Languages used on the internet today

---

5- Languages, <https://www.internetworldstats.com/stats7.htm> , 2018.

# Colloquial Arabic language

Analyzing Arabic contents is a **challenging task** due to:

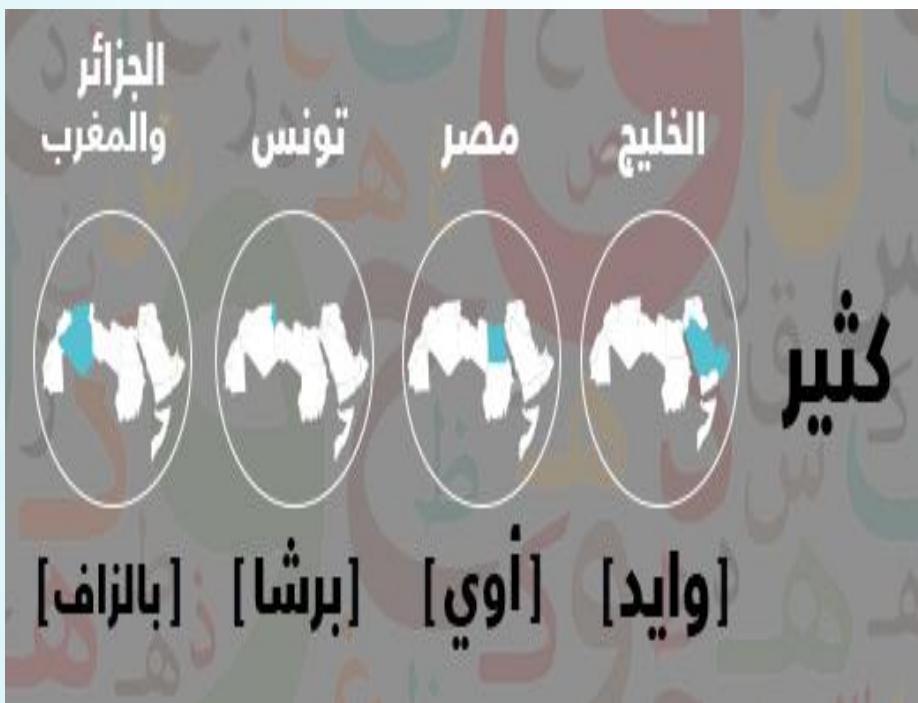
- Unstructured language (informal language).
  - Spelling inconsistencies.
  - Slang words. Such as (يلا, حلو، حلوس)
  - Idiomatic expressions.

بعد ما شابه ودوه الكتابة، التكرار يعلم الممار، الطيور على اشكالها تقع ( طبائع السم بيدوقة، شعارات ونرزمي، هاته الشيل من ذلك الأسد

# Colloquial Arabic language

- “Colloquial Arabic” (CA) is most widely spoken and written in daily live conversations with different dialects.

( colloquial Arabic varies from region to another region, from country to other countries, even from even state to another.)

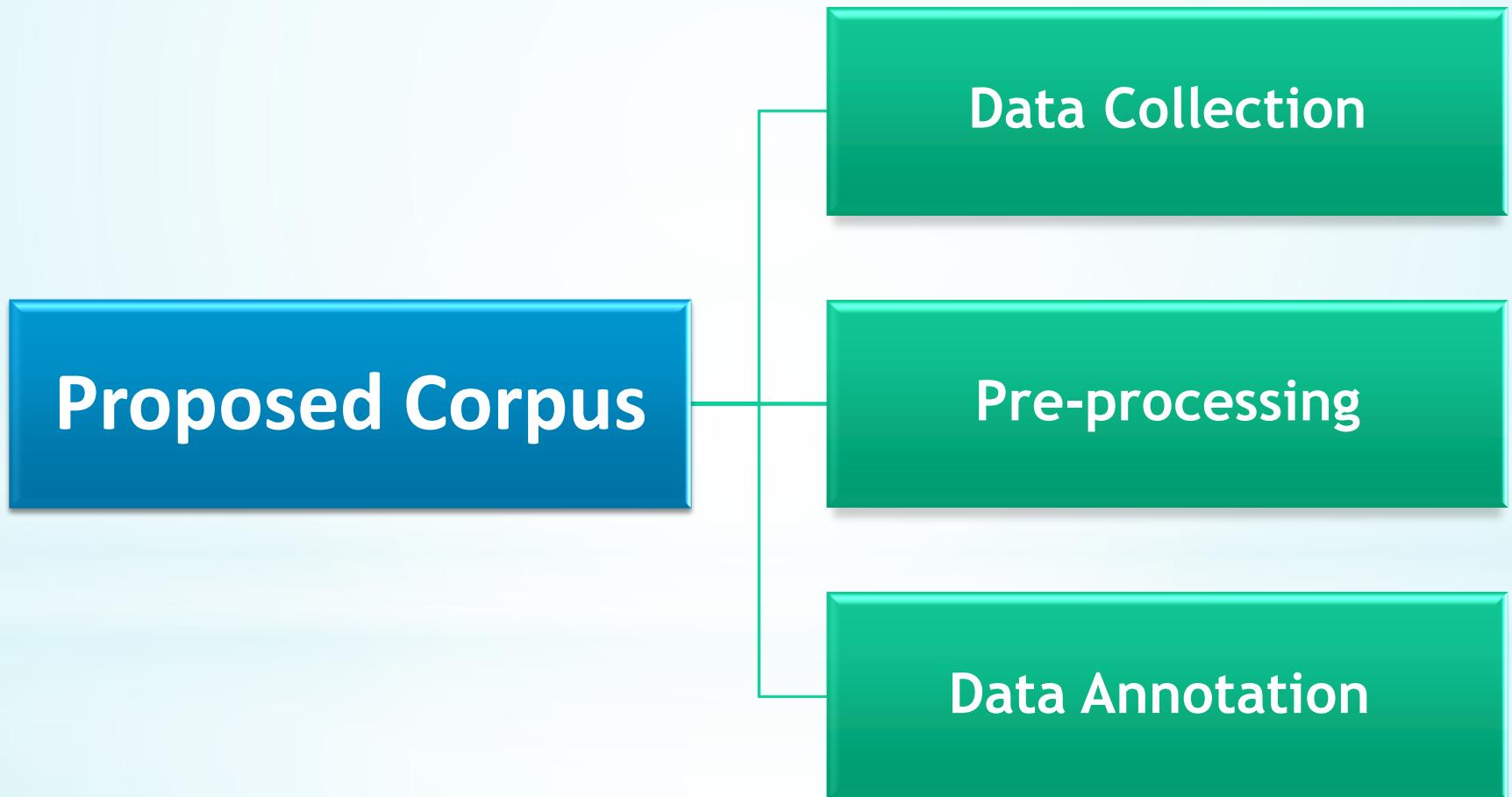


Colloquial dialect	Sentence
Jordanian	أيش بدهك
Egyptian	عايز ايه
Saudi	وش تبغى
Tunisian	شنو تحب
Algerian	واش تحب

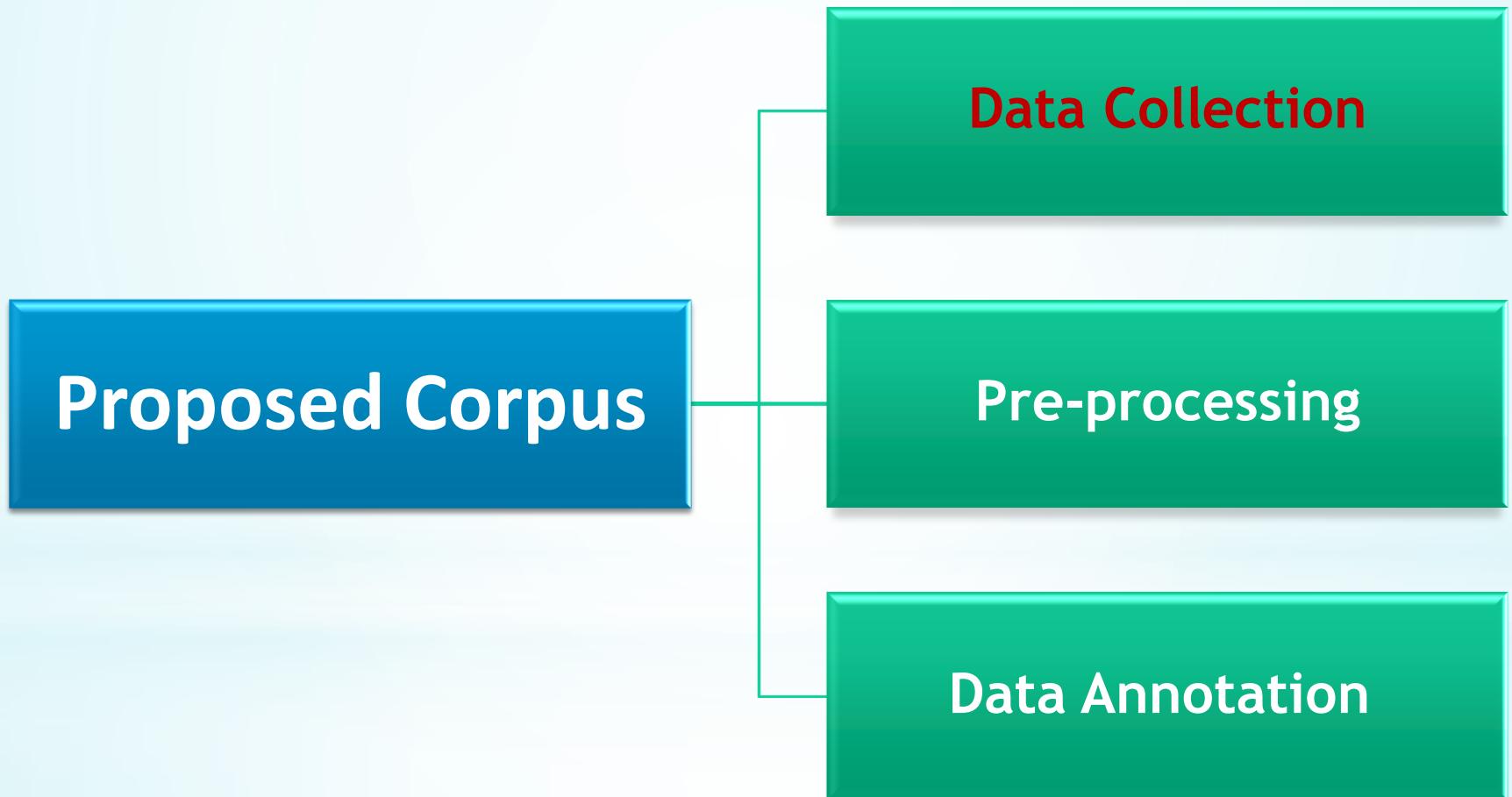
# Agenda

1. Introduction.
2. **Proposed Corpus.**
3. Embedding Layer.
4. Architecture of RCNN Model.
5. Publications.
6. RCNN Code.

# Proposed Corpus (Egyptian Corpus)



# Proposed Corpus (Egyptian Corpus)



# Proposed Corpus (Egyptian Corpus)

- ✓ Tweets was collected using twitter APIs.
- ✓ The collected tweets span the period from April 11, 2015 to December 12, 2015. In this initial stage, a collection of 2,408,128 noisy raw tweets were gathered.

Mohamed Salah @MoSalah

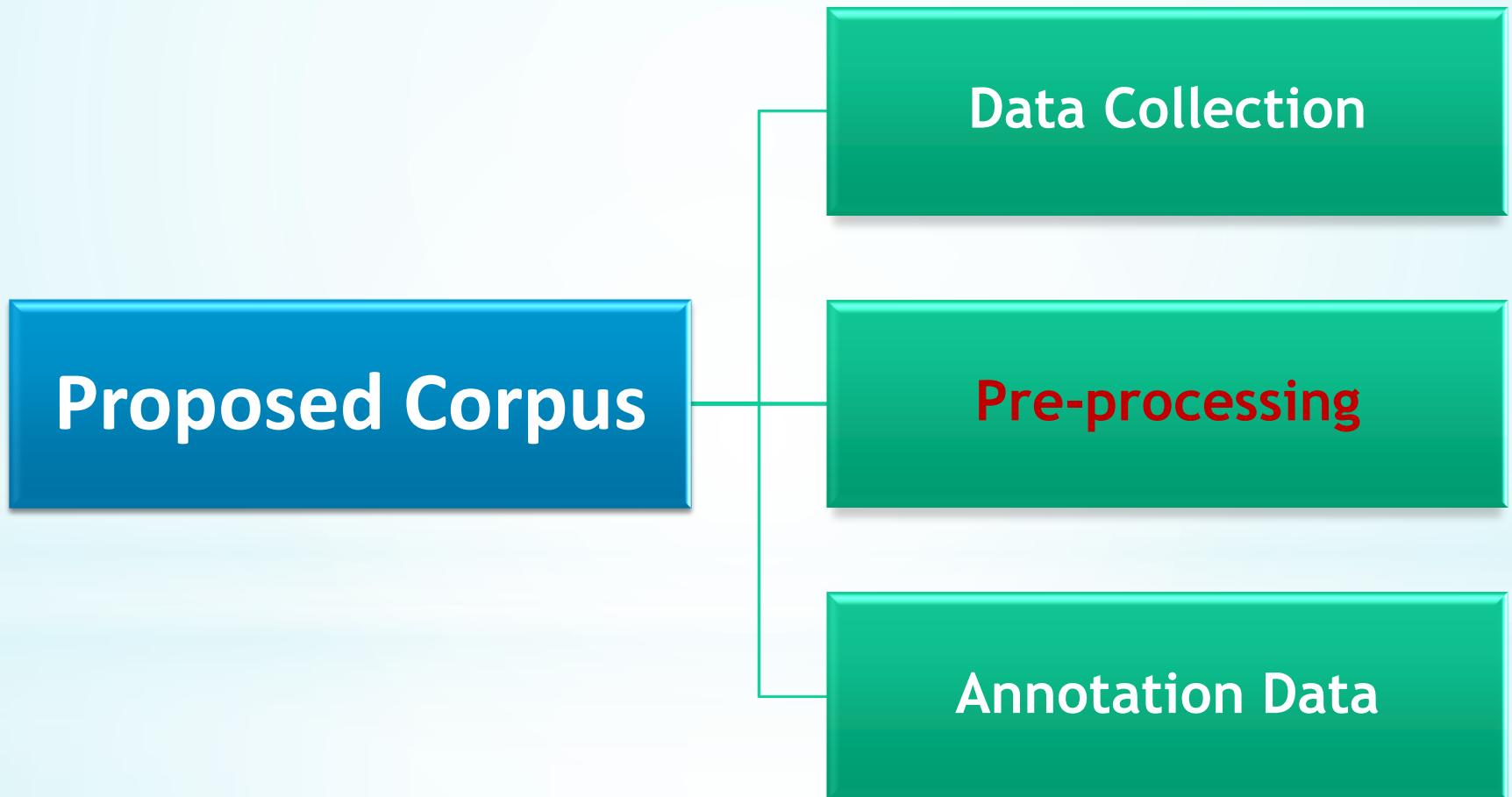
مهما كانت الظروف يظل الأهلي النادي الأكبر في أفريقيا والوطن العربي .. الأسد يمرض ولا يموت.

5:44 am · 28 Nov 2018 · Twitter for iPhone

5.3K Retweets 95.9K Likes

## Original tweets

# **Proposed Corpus (Egyptian Corpus)**

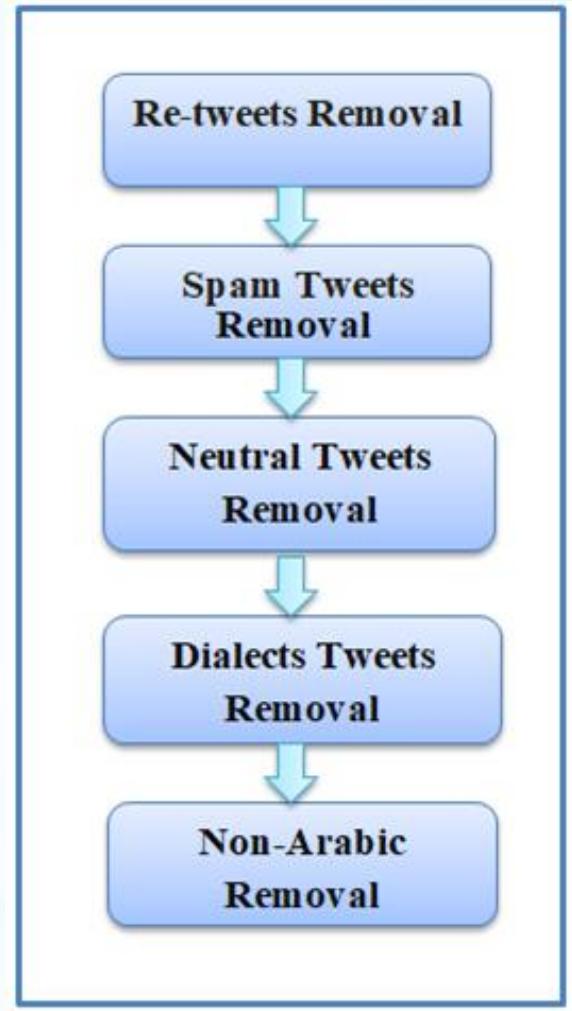


# Preprocessing (Egyptian Corpus)

- Phase 1: Filtering Data
- Phase 2: Cleaning Data
- Phase 3: Normalization Data

# Phase 1: Filtering Data

- After collecting the data from Twitter, we found that the collected tweets have **many problems** that will affect the annotation process accuracy.
- This phase consisted of **five** steps **manually**:



**Tweet filtering phase**

# Phase 2: Cleaning Data

- ✓ This phase is important since it deals with the **noisy** nature of tweets.
- ✓ data cleaning deals manually:

Such as (hashtag (#), username (@username), emoticons (=D, XD), time stamps and URL)

- ✓ Some tweets may contain:  
(missing spelling words, missing letters and written incorrectly)
- ✓ This **problem** is handled by **correcting** these **words** and wrong characters.
- ✓ Some words are **combined together** due to space limitation of Twitter. This problem is handled by **adding spaces** between words.

# Phase 3: Normalization Data

- Normalization is the process of transforming the Arabic text in order to be **consistent**.

For example:

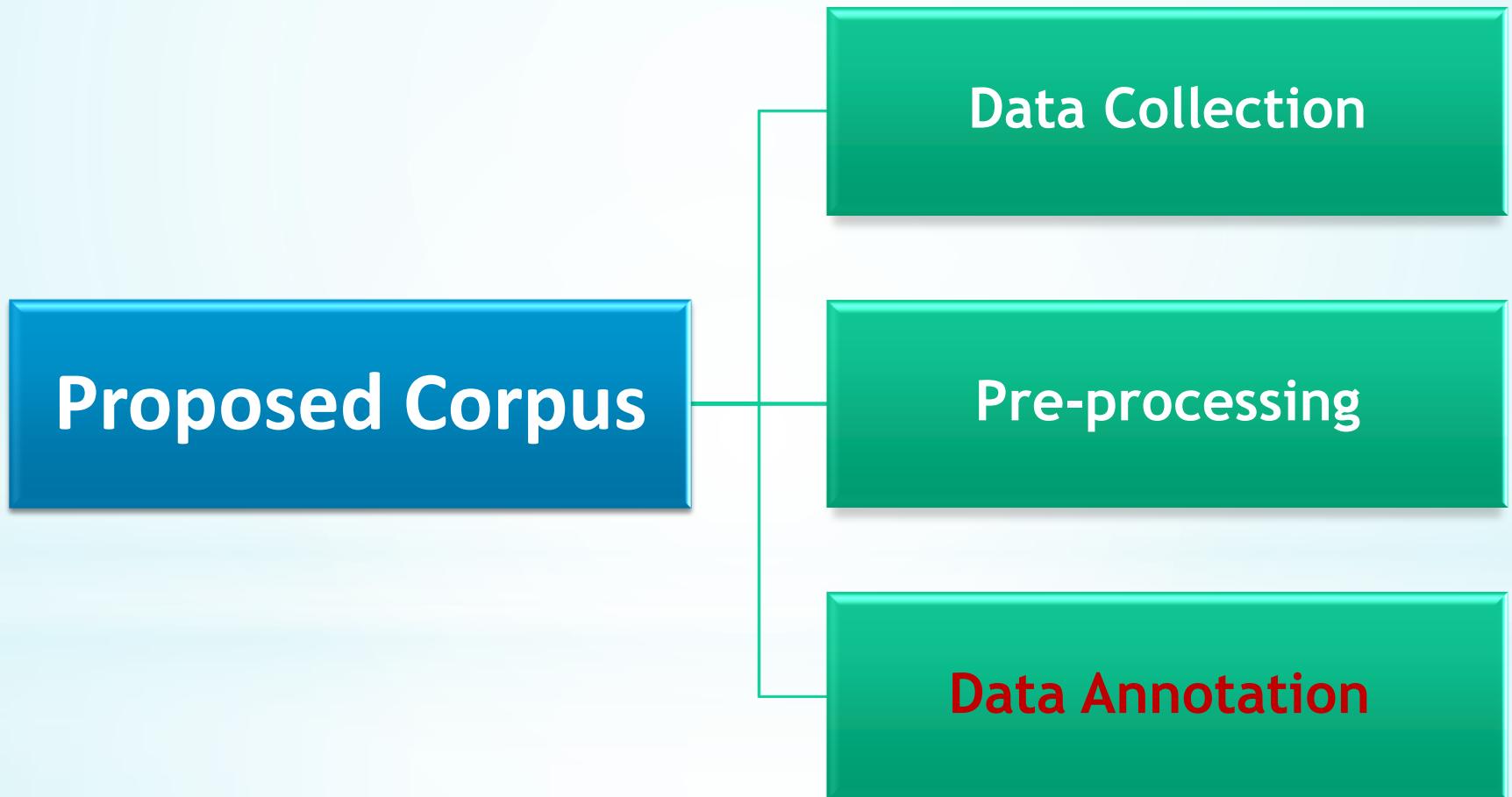
انت	إنت	انت	أنت
-----	-----	-----	-----

All these forms cause the word to be considered as three different words.

## Normalization Rules

Rule	Example
Tashkeel	> تُقْرَبُ تقرب
Repeate	>-كَتِيررررر كتير
Tatweel	->اللَّهُ الله
Hamza	(ء ، ي ، و ) -> ء
Alef	(ا ، اً) -> ا
Lamalef	(ع ، لـ) -> لا
Yeh	ي -> (ى ، ي)
Heh	(ؤ ، ة) -> ئ

# Proposed Corpus (Egyptian Corpus)



# Data Annotation

- At this phase, the annotation process is done **manually**, which manually classify each tweet as “*positive*” or “*negative*”.



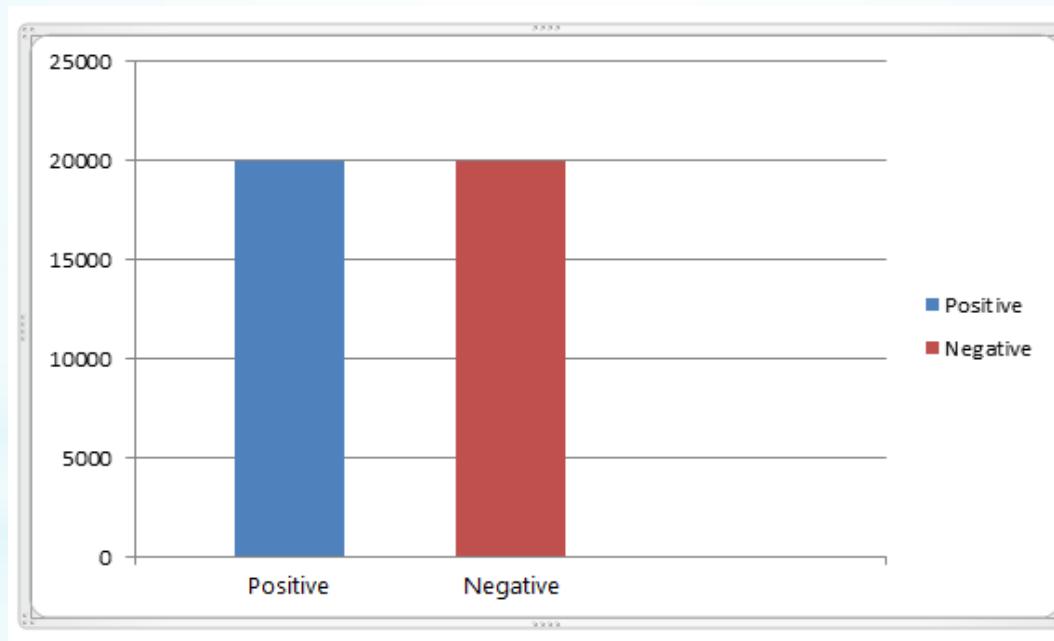
# Egyptian Corpus

# Proposed Corpus (Egyptian Corpus)

- ✓ New Arabic corpus for mixed tweets written in Egyptian dialect and MSA.
- ✓ Cover **several topics including:**
  1. Political, Sportive, Health, Social topics.
  2. Sarcastic jokes.
  3. Proverbs.
  4. Poetry.
  5. Users' opinions concerning different products.

# Proposed Corpus (Egyptian Corpus)

- ✓ This corpus consists of 40000 tweets (classified as 20000 positive and 20000 negative).



Total distribution of Corpus (balanced data)

# Proposed Corpus (Egyptian Corpus)

✓ It divided into two columns:

- The first column is called **review**.
- The second column is called **label**.

	B	A	
	negative	ده اللي انت فهمته . ياريتها كانت بركت عليك وهي بتولك .	18148
	negative	بلاش الهجوم الشرش ده . علشان حد اختف معاك في الراي .	18149
	negative	نربتها اللي ما يموش من الاصاص يومت مثلول .	18150
	positive	اذا كانت تقرأ فاني اذكر فلن ان و اذا كانت ايشست فاني احي بابتسامتك هذه .	18151
	negative	في كمية تعلم و غباء و جهل حاجة ترفع المضetto .	18152
	negative	اهم اشخاص في حياتي . هم ابعد ناس عنى .	18153
	negative	با خرابي على الدماغ .	18154
	negative	انا بيت عنان مردتش عليك .	18155
	negative	هي افراح العالة مش بتخل على الا أيام الامتحانات .	18156
	negative	للرجل حكليات بعضها انختارها وبعضاً نجر علىها .	18157
	positive	ابدي دي هي اللي هنجيب الملايين بعد كده .	18158
	negative	المواقف بتينين مين اللي هيسترجل في وقت الشدة و مين بيتمثل .	18159
	positive	الله يخليكي ، ويجعلها سنة سعيدة منهاش صداع في راسنا بارب .	18160
	positive	الحمد لله والله انتي نعمة من ربى .	18161
	positive	بصراحه انسانه جده جداً و رقيقة جداً و محترمة جداً و خدومه . انا اشرفت اني اتعرفت عليك .	18162
	negative	كل بنت غير ايه مني صدقني حقاً .	18163
	positive	انا احتاج شخص يخليني احب الحياة من جديد .	18164
	positive	بارب بجي اليوم الى البس فيه كا .	18165
	positive	ده كلام لوجيك جدا .	18166
	negative	عليزه اكتب توينه عميقه بس دماغي واقفة .	18167
	positive	رينا يسعدكم و يرزقكم بصبيان و بنات .	18168
	positive	بارب انزل على قلبى حب المذكرة .	18169
	negative	الاسعار كل يوم تزيد . حرام بجد .	18170
	positive	بعندي تخيلوا حد يخطب عليك ويقولك هابي نبو بير و بديك هديه . معالترشح ببابا نويل للرئاسه .	18171
	positive	انفسنلوا معانا احتنى .	18172
	negative	لولا الخيل كان مونتنا من الواقع .	18173

## Sample of our corpus

D	C	B	A	review
Experts#2	Experts #1	label		1
N	N	negative	اكبر خطأ ترتكبه ان تعامل الناس باخلاقك انت مش باخلاقهم هما .	2
N	N	negative	دائما اكره اخر ليله في كل مكان .	3
N	N	negative	يارب اللي يسرق تويتاتي يدخل النار .	4
N	N	negative	الاسف في تناول القهوة يسبب الوفاة .	5
P	P	positive	انا ابيهيلك من التراب النهارده حاجة تعرف .	6
N	N	negative	في بنات بتلبس اكسسوارات لندرجه انك تحس ان في حنطور ماشي جنبك .	7
P	P	positive	احتاج صديق حقيقي يواسيبني ويخفف عنى .	8
P	P	positive	لازم اتعلم الثبات الانفعالي زيه كدا .	9
N	N	negative	جروب الفعه ليلة الامتحان دايما تحس ان الناس اللي فيه بيتذكرة منهجه تاني . وبيكتهو تعاوين على الجروب .	10
N	N	negative	النقاش هو الذي يؤدي الى نتيجه اما التصعيد ففيولد تصعيد اكتر منه .	11
N	N	negative	وصلت الى مرحلة اني كرهت بلدي واتي فيها .	12
P	P	positive	عندما فرأت قصه اول شهيد بالإسلام ام عمار عرفت كم نحن ضعاف .	13
P	P	positive	اجمل انسان هو المتسابح ، المتواضع .	14
N	N	negative	اسمها شبكة تواصل اجتماعيه . مثل شبكه مراقبه اجتماعيه .	15
N	N	negative	امي عندها صور لي تنهى مستحبى .	16
N	N	negative	دنيا مليئان أمان صحيح .	17
P	P	positive	الرئيس الخدين رزق .	18
N	N	negative	الحكم اللي انتقال النهارده عليهم قاسي جدا .	19
N	N	negative	عند الاستيقاظ تكتب لمهامك الجديدة .	20
N	N	negative	مصر في المرتبه السابجه في ترتيب الدول الأكثر خطورة في العمل الصناعي انا ايه اللي دخلني اعلام .	21
N	N	negative	سيناريو عقري وخطة محكمة شقيق ينجح في نفس اليوم اللي مبارك يموت فيه . وشقيق يتحوز سوزان وجمال بورث الحكم .	22
N	N	negative	العرب يسعون لنفس الاحلام الغرب يسعون لتحقيق الاحلام هذا هو الفرق بين العرب والغرب .	23
N	N	negative	حقيقة غاب صدام ظهر الشيعة .	24
N	N	negative	بعد هذه الاختيارة لدعم مرسي وهذا الصوت المبهر . اؤكد لكم ان احمد شقيق رئيسا للجمهورية باكتساح .	25
N	N	negative	قد تكون حياتي مليئة بالأصدقاء ولكن قليل ما يشعر بي احد .	26

## Notes:

- ❖ we construct the corpus from those tweets having clearly positive or negative sentiment.
- ❖ Also, to validate the annotation of corpus, two different experts were asked to check the annotation. The results of both experts were consistent with the corpus annotations with 100 % accuracy.

# Egyptian Corpus Statistics

Total number of tweets	40000
Number of positive tweets	20000
Number of negative tweets	20000
Number of words	359,818
Max tweet token	39
Number of tokens	1,953,869
Average tokens per tweet	17

✓ This corpus is available in [6].

Harvard Dataverse > Corpus on Arabic Egyptian tweets

 Info – The "DRAFT" version was not found. This is version "1.0".

## Corpus on Arabic Egyptian tweets

Version 1.0

Rania Kora; Ammar Mohammed, 2019, "Corpus on Arabic Egyptian tweets", <https://doi.org/10.7910/DVN/LBXV9O>, Harvard Dataverse, V1

 Cite Dataset ▾

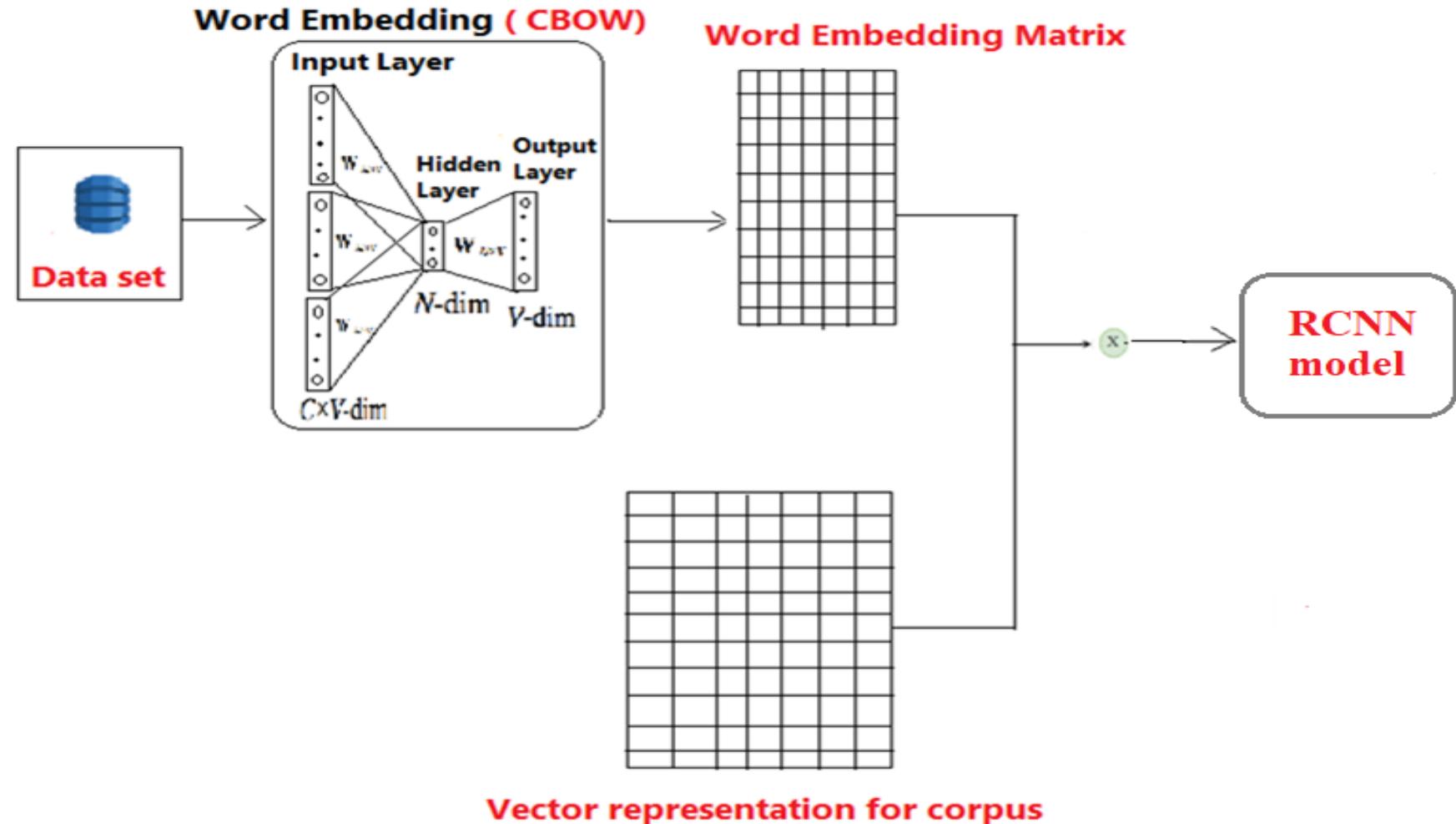
Learn about Data Citation Standards.

- 
- 6- Kora, R. and A. Mohammed, “Arabic tweets Egyptian dialect”, . Mendeley Data, V1, 2019.  
<https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LBXV9O&version=DRAFT&faces-redirect=true>

# Agenda

1. Introduction.
2. Proposed Corpus.
3. Embedding Layer.
4. Architecture of RCNN Model.
5. Results and Discussion.
6. Publications.
7. RCNN Code.

# Embedding Layer of RCNN Model



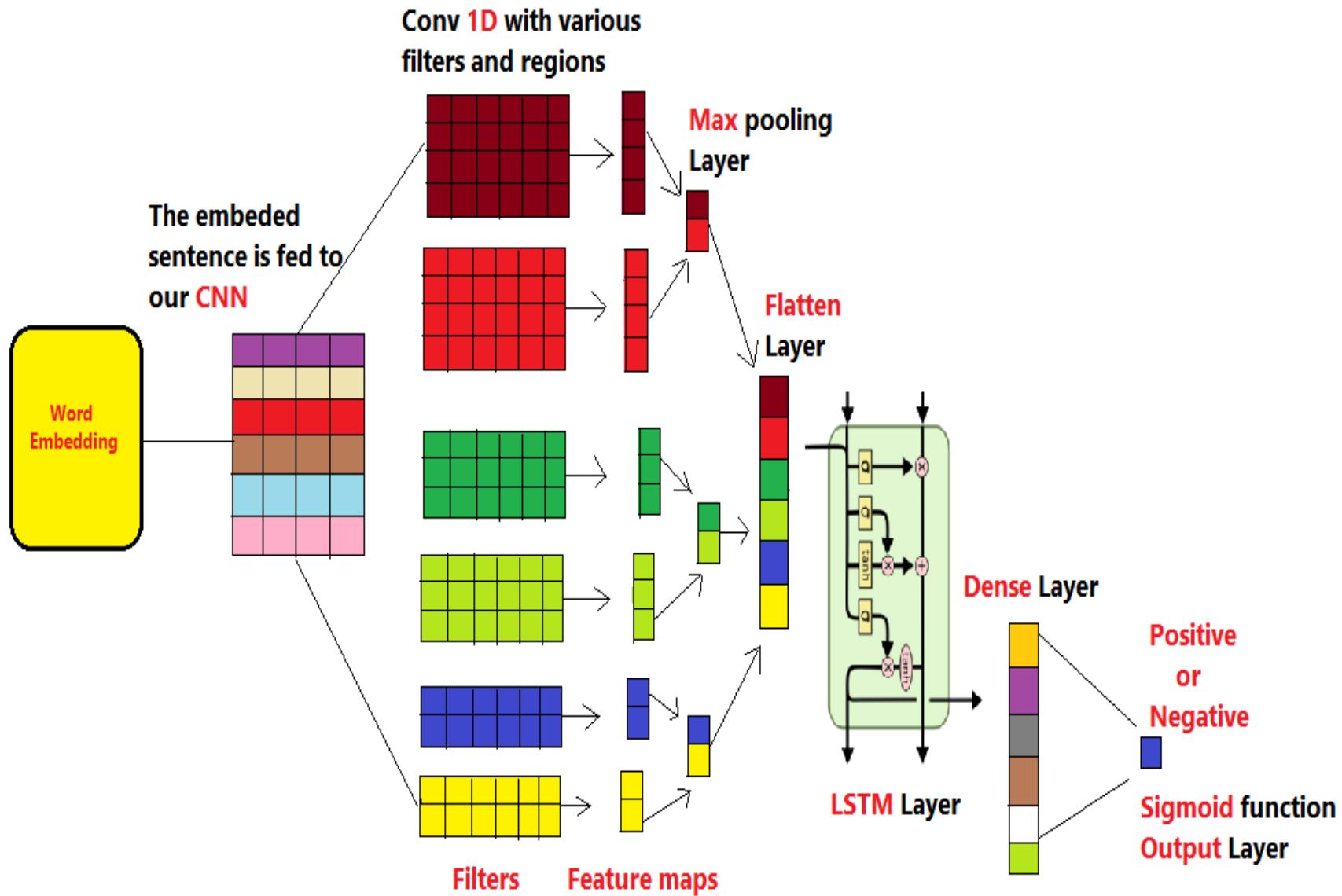
# Agenda

1. Introduction.
2. Proposed Corpus.
3. Embedding Layer.
4. Architecture of RCNN Model.
5. RCNN Code.
6. Publications.

# RCNN Model

- Recurrent Convolutional Neural Network (**RCNN**) is a network combining the advantages of both **CNN** and **LSTM** neural networks.
- The RCNN architecture consists of: Embedding Layer, Convolutional Layer, Max-pooling Layer, Flatten Layer, LSTM layer, and Dense an activation layer.

# Architecture of RCNN Model



# Agenda

1. Introduction.
2. Objectives.
3. Proposed Corpus.
4. Proposed Approaches.
5. Publications.
6. RCNN Code.

# Agenda

1. Introduction.
2. Proposed Corpus.
3. Embedding Layer.
4. Architecture of RCNN Model.
5. Publications.
6. RCNN Code.

# Publications

1. Mohammed, A., Kora, R.: *Deep learning approaches for arabic sentiment analysis.* *Social Network Analysis and Mining* 9 (1), 52 (2019) Springer. DOI [10.1007/s13278-019-0596-4](https://doi.org/10.1007/s13278-019-0596-4). URL <https://link.springer.com/article/10.1007%2Fs13278-019-0596-4>

Social Network Analysis and Mining (2019) 9:52  
<https://doi.org/10.1007/s13278-019-0596-4>

ORIGINAL ARTICLE



## Deep learning approaches for Arabic sentiment analysis

Ammar Mohammed<sup>1</sup> · Rania Kora<sup>1</sup>

Received: 27 June 2019 / Revised: 3 August 2019 / Accepted: 9 September 2019  
© Springer-Verlag GmbH Austria, part of Springer Nature 2019

### Abstract

Social media are considered an excellent source of information and can provide opinions, thoughts and insights toward various important topics. Sentiment analysis becomes a hot topic in research due to its importance in making decisions based on opinions derived from analyzing the user's contents on social media. Although the Arabic language is one of the widely spoken languages used for content sharing across the social media, the sentiment analysis on Arabic contents is limited due to several challenges including the morphological structures of the language, the varieties of dialects and the lack of the appropriate corpora. Hence, the rapid increase in research in Arabic sentiment analysis is grown slowly in contrast to other languages such as English. The contribution of this paper is twofold: First, we introduce a corpus of forty thousand labeled Arabic tweets spanning several topics. Second, we present three deep learning models, namely CNN, LSTM and RCNN,

Activate Windows  
Go to Settings to activate Wi

**2. Kora, R., Mohammed, A.: *Corpus on Arabic Egyptian tweets* (2019). DOI 10.7910/DVN/LBXV9O. URL <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LBXV9O&version=DRAFT&faces-redirect=true>**

 HARVARD  
Dataverse

Search ▾ About Use

Harvard Dataverse > **Corpus on Arabic Egyptian tweets**

 Info – The "DRAFT" version was not found. This is version "1.0".

 **Corpus on Arabic Egyptian tweets**  
Version 1.0

Rania Kora; Ammar Mohammed, 2019, "Corpus on Arabic Egyptian tweets", <https://doi.org/10.7910/DVN/LBXV9O>, Harvard Dataverse, V1

 Cite Dataset ▾ Learn about Data Citation Standards.

# Agenda

1. Introduction.
2. Proposed Corpus.
3. Embedding Layer.
4. Architecture of RCNN Model.
5. Publications.
6. RCNN Code.

thank you

