

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.DOI

# Social Media and Microblogs Credibility: Identification, Theory Driven Framework, and Recommendation

KHUBAIB AHMED QURESHI<sup>1</sup>, RAUF AHMED SHAMS MALICK<sup>2</sup>, MUHAMMAD SABIH<sup>3</sup>,

<sup>1</sup>DHA Suffa University, Karachi, Pakistan (e-mail: k.ahmed@dsu.edu.pk)

<sup>2</sup>National University of Computer & Emerging Sciences, Karachi, Pakistan (e-mail: rauf.malick@nu.edu.pk)

<sup>3</sup>DHA Suffa University, Karachi, Pakistan (e-mail: m.sabih@dsu.edu.pk)

Corresponding author: Khubaib Ahmed Qureshi (e-mail: k.ahmed@dsu.edu.pk).

**ABSTRACT** Social media microblogs are extensively used to get news and other information. It brings the real challenge to distinguish that what particular information is credible. Especially when user authenticity is hidden, due to the microblog's anonymity feature. Low credibility content creates an imbalance in society. Therefore many research studies are conducted to assess automatic microblog's credibility but the majority of them offer different concepts of credibility and the problem seems unresolved. Credibility is multi-disciplinary, hence there is no generalized or accepted credibility concept with all its necessary and detailed constructs/components. Therefore, it is necessary to understand the complete anatomy of information credibility from different disciplines. It is accomplished here through an in-depth and organized study of all the problem dimensions for the identification of comprehensive and necessary credibility constructs. The framework is also proposed based on the identified constructs. It adheres to these constructs and presents their inter-relationships. It is believed that the framework would provide the necessary building blocks for implementing an effective automatic credibility assessment system. The framework is generic to social media and specifically implemented for microblogs. It is completely transformed up to features level, in the context of microblogs. Regarding automatic credibility assessment, it is proposed after detailed analysis that the attempt should be made for hybrid models combining feature-based and graph-based approaches. It is observed that quite a few surveys in the literature focus on some limited aspects of microblogs credibility but no literature survey and fundamental study exists that consolidates the work done. To understand the broader domain of credibility and consolidate the work in this area that can lead us to a suitable framework, we explored the existing literature from different disciplines for the said objectives. We categorized them along various dimensions, developed taxonomy, identified gaps and challenges, proposed a solution, developed a theory-driven framework with its transformation to microblogs, and suggested key areas of research.

**INDEX TERMS** Social Media Credibility, Twitter Information Credibility, Credibility Features, Automatic Credibility Assessment Models, Proposed Solution, Credibility Framework, Credibility Taxonomy, Credibility Levels Dimensions Constructs, Credibility Studies, Credibility Dataset.

## I. INTRODUCTION

MICROBLOGS are intensively used to share news [1], opinions, observations, health issues, entertainment, experiences, and many more [2]. It is therefore becoming an imperative source of information but on the other hand, not-credible [3], [4] and cumbersome [5]. Taking an example of microblogs such as Twitter is steadily achieving gigantic consideration [6] as an important form of information media [7]. A large number of users throughout the world spread a wide range of information in real time [8]. Millions of Tweets are posted per hour on Twitter. Currently, it is the growing social medium and prevalent news media source as well [9]. Users massively share news headlines and also report real-time events of varying nature, well before official sources [8]. Twitter users are of many kinds, such as citizens, companies, governments, famous personalities, politicians, and many more, and such a wide range of users heavily depend on it for their business, political, social, and educational communications. Therefore on the dark side of this beautiful picture spammer also exploits the anonymity feature of microblogs

to propagate their spam messages and scam URLs. It is quite vulnerable and turns into a medium of wrongdoers to spread rumors, fake news and other forms of misinformation [10]–[13]. Spread of hate speech [14], [15], political astroturf memes [16], extreme biases [17] are also found. Low credibility content creates an imbalance in society by damaging the reputation, public trust, freedom of expression, journalism, justice, truth, and democracy. Consequently, microblogs' users often need to judge the information's credibility. It becomes more challenging when source/user authenticity is hidden from the viewer, though user anonymity is one of the prose of microblogs. Unfortunately, it also welcome some other issues like: user's coordinated behavior [18], follower's fallacy [19], etc. It not only affects the quality of microblogs content but also introduces another challenge for gauging the source credibility.

There are many studies conducted at different aspects of credibility in many fields, such as; psychological factors affecting credibility, credibility types, dimensions, constructs, theoretical credibility frameworks, user's perceptions of credibility, suggested credibility features, automatic credibility assessment studies, and experimental studies of ranking information based on credibility, etc. Even then there is neither comprehensive nor accepted credibility attempt exist [20], [21], nor there is a standard definition of credibility found, though there are some related terms used to define credibility [22]. Considering the broader domain of credibility, having related terms or even having definitions only, never provides us that these are the necessary aspects that must be considered when credibility is assessed. Though it is required and extremely important in doing such assessments. In continuation with these challenges. It is also discovered that no literature survey and fundamental study exists that consolidates the work done from different fields. Therefore to fulfill the objectives. The literature is explored to identify such necessary credibility components. These identified components also lead us to propose a suitable framework of automatic credibility assessment.

Another very obvious fact to be highlighted to understand the importance and need of such broad and in-depth study; is about different types of malicious profiles or simply called malicious accounts. Which are completely ignored in all credibility studies. Though there are separate bot-detection studies found but not under the umbrella of credibility or not considered as a necessary aspect of credibility. Examples of such malicious profiles are; Bots, Trolls, Cyborgs, etc. All such forms of malicious profiles are usually believed to aggravate the wrong sense of credibility indicators and play a key role in the spread of low credibility contents [23]. It became very evident in investigations into Russian attempts to influence the 2016 US election [24]. It has also been observed that a massive amount of low credibility contents have already been shared over social media and microblogs before and after the US Election 2016 despite many efforts of credibility assessments [25]–[28]. It shows that some important and necessary aspects were ignored in available

credibility assessment methods, as discussed earlier.

Although credibility has been studied since ancient times,



FIGURE 1. Majority of the studies only cover either one or only some of the above aspects of credibility and a majority of the aspects are left undiscovered.



FIGURE 2. Above are some general aspects of credibility which are completely missed in literature within the context of credibility. Low credibility contents may have the above forms, which should also be considered when credibility assessment is made.

and in different research fields to date, such as psychology, media science, information science, communication, journalism, social sciences, and information retrieval, etc. [29]. It is noticed in literature that, due to being multi-perspective nature the diversity in the definition and perception of credibility reflects different viewpoints in different work studies. These studies only stick to just a single or only a few aspects of credibility. Some studies consider only Relevance as a criterion of being credible, some assume just Reputation as the major driver of Credibility, whereas the majority only stick that Fake and Rumor identification is credibility identification. It is also perceived by researchers, that Rankings concerning author Influence and Topic Expertise are strongly treated as credibility ranking. The majority of studies exploit just Informativeness as a credibility indicator. Few found examining Trust level as true credibility judgment. It is observed and quite evident in many research studies as well, that the credibility notion needs to be standardized because many studies only cover either one or some aspects of credibility (see figure 1) and a majority are left undiscovered. Some potential aspects are not even explored though much affect the credibility (see figure 2). Effective and comprehensive credibility concept may conforms some combined aspects

presented in both figure 1 and figure 2. It means that low credibility contents may have a variety of forms presented in both figures. There is another strong observation developed through a majority of research studies, that credibility is assessed for news contents only (fake/real), though it equally exists in non-news contents as well, with a different set of aspects. Therefore those necessary set of credibility related aspects need to be identified which must be evaluated for any piece of information in terms of its credibility assessment. It is already discussed that credibility is multi-disciplinary, hence there is no generalized or accepted credibility concept with all its necessary and detailed constructs/ components. It is extremely necessary and quite challenging, to understand the broad domain of information credibility to extract its complete anatomy from different disciplines. It could be accomplished through an in-depth and organized study of all the problem dimensions and identification of comprehensive and necessary credibility constructs under credibility's definition first. Further, the development of a concrete framework that adheres to those basic constructs/components could be possible. The framework will be theory-driven and provide a complete relationship/connection between different identified credibility components. In this study, we are concerned with the said identification followed by the development of a generic and comprehensive framework of information credibility. The framework will be generic to social media and specifically implemented for microblogs. It will be completely transformed up to features level, in the context of microblogs.

Nowadays numerous applications use a vast amount of microblogs data, such as; recommendation systems, event detection systems, social bookmarking systems, disaster response applications, campaign management systems, business monitoring applications, different types of prediction systems, and microblog search engines, etc. Each one of them only requires credible data to make these systems more effective [30]–[32]. Therefore dealing with information credibility problems in microblogs and social platforms, is necessary [33]. Once we would be able to develop an efficient and comprehensive credibility framework, which is missing and required, then there could be many applications in which the credibility framework would successfully contribute. For example; one of the most obvious applications could be the determination of the credibility of various posts during major global or local events. This can help for example in disaster response situations where the important information such as the extent of damage and need for action, can be figured out based on a large amount of microblogs posts and the trust ratings of their posters.

It is observed that quite a few surveys in the literature focus on some limited and individual aspects of microblogs credibility like health info. credibility [29], user influence/source credibility [34], trust in social networks [35], relevance-trust and influence [36]. There is a surface level or extremely short survey conducted over twitter information credibility in [37] and another general survey over information credibility of

social media is done in [38]. As far as we discovered that there is no literature survey and fundamental study exists that consolidates the work on credibility similar to this study.

The remaining of the paper is organized as follows: problem formulation is done in section II. Credibility Taxonomy is developed as table:1 and figure: 3. The same is discussed from section 3-7, such as: in section III different definitions of credibility with its necessary and related components (levels, dimensions, and constructs, etc.) are presented. It helps us to understand credibility in the broader sense. Section IV highlights theoretical credibility frameworks. The most important section V presents many research areas which must be considered in credibility study and found extremely supportive, therefore named as supported research. Taxonomy's main section VI purely focuses only on all social media and microblogs specific information credibility studies. Last section VII of taxonomy is about standard credibility datasets. Section VIII literature-based important features are presented. Section IX summarizes the study through important findings and discussions. In section X we presented first, all theories in support of credibility framework identification and then our proposed theory-driven credibility framework is presented in section XI followed by section XII as Recommendations. Section XIII is about future research directions and section XIV concludes our study. Challenges and limitations are presented within different sections. Important terms used in the study are defined in appendix.

## II. PROBLEM FORMULATION:

To better understand the problem, in this section, we have formulated the credibility assessment as a classification problem and scoring/ranking problem. The mathematical problem formulation is done as following:

Let  $P = \{p_1, p_2, \dots, p_n\}$  be the set of  $n$  Posts, and  $U = \{u_1, u_2, \dots, u_m\}$  be the set of  $m$  Users on microblog. Each  $p_i$  consists of series of features including text domain, text sentiment score, text length, post spread score, no. of comments and replies, etc. Similarly each  $u_i$  consists of series of features like: influence score, name, domain, date creation, etc.

Classification Problem: Given Post  $P$ , and User  $U$  goal is to learn prediction function, such as  $f(p_i, u_j) \rightarrow \{0, 1\}$  satisfying:

$$f(p_i, u_j) = \begin{cases} 1 & \text{if } p \text{ is credible} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Scoring/Ranking Problem: This could also be ranking/ scoring function, such as:  $f(p_i, u_j) \rightarrow \{0, 1, 2, 3, 4, 5\}$  satisfying:

$$f(p_i, u_j) = \begin{cases} 0 & \text{if } p \text{ is not - credible} \\ 1 & \text{if } p \text{ is low - credible} \\ \cdot & \cdot \\ \cdot & \cdot \\ 5 & \text{if } p \text{ is highly - credible} \end{cases} \quad (2)$$

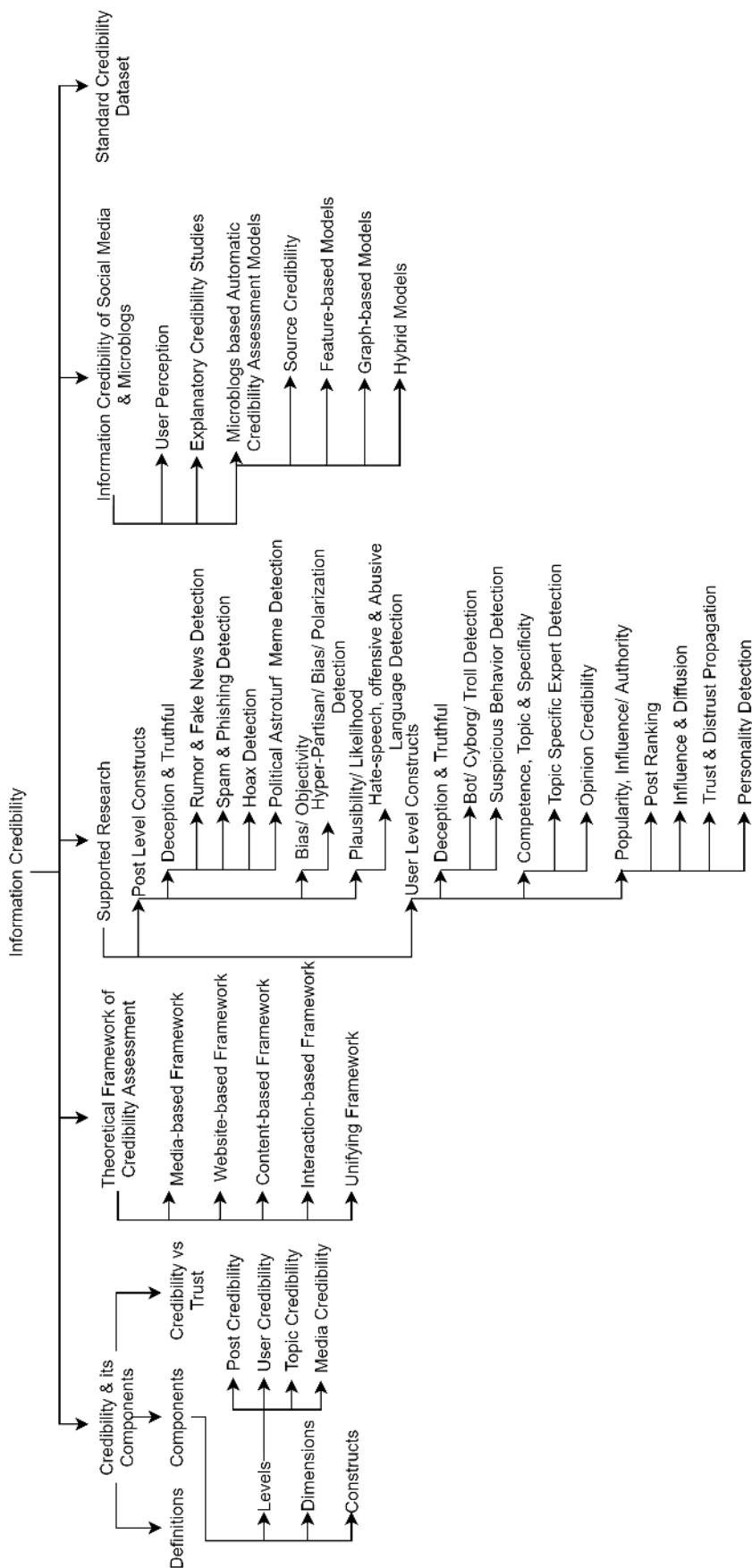


FIGURE 3. Detailed Credibility Taxonomy: the organized and complete taxonomy with all its levels is presented in this figure.

**TABLE 1.** Simplified Credibility Taxonomy: only top level and lowest levels are presented in this tabular form, intermediate levels are explicitly omitted for simplicity and better understanding. Detailed taxonomy with complete levels are shown in credibility taxonomy figure 3.

S. No.	Category	Sub-Category	Reference	Description
1	Credibility Definitions	Believability, Trust, Reliability, Accuracy, Fairness, Objectivity	[39]	How credibility is defined and its related components, e.g.: Levels, Dimensions, & Constructs, etc. and what is the relationship between credibility and trust.
		Quality of being Trusted and Believed	[40]	
		Quality of being Believed	[41]	
		Credibility has components: Message, Source and Media	[42]	
		Expertise and Trustworthiness	[43]–[46]	
	Believable Person and information	[47]		
	Credibility Components	Levels, Dimensions, and Constructs	see table 2, 3	
	Credibility VS Trust	Credibility is antecedent to Trust	[48]–[51]	
2	Theoretical Credibility Assessment Framework	Media-based Framework	[52]–[55]	These conceptual or theoretical frameworks provides: 1. Categorization similar to evolutionary generations. 2. Understanding of credibility assessment process & related concepts & how it is affected. 3. Underlying process involved behind people to perform assessment of credibility.
		Website-based Framework: Fogg’s Prominence Interpretation Theory	[43]	
		Content-based Framework	[56], [57]	
		Interaction-based Framework: Rieh’s Predictive and Evaluative Judgment	[58]	
		Interaction-based Framework: Wathen, Burkell- First Medium is rated, then source and message, third interaction of presentation and content	[59]	
		Interaction-based Framework: Sundar’s MAIN model (Modality, Agency, Interactivity, and Navigability) four “affordances” in digital media	[60]	
		Interaction-based Framework: Elaboration Likelihood Model (ELM) of Persuasion	[61]	
		Interaction-based Framework: Heuristic Systematic Model (HSM) of information processing	[62]	
		Interaction-based Framework: Controlled and Automatic Processing Models (CAPM)	[63]	
		Interaction-based Framework: Social Information Processing Theory (SIPT)	[64], [65]	
		Interaction-based Framework: Dual processing model for Web	[66]	
		Unifying Framework: Provides basic levels: Interaction, Heuristics, and Construct	[67]	
		Unifying Framework: Rieh et al- Extension	[68]	
3	Supported Work	Misinformation/Disinformation: Rumor and Fake News Detection	[3], [69]–[76] and [25], [77]–[85]	These are all studied as separate research areas in the literature, though each one of them are different construct/ aspect of credibility, therefore we consider them as important building blocks of credibility or necessary components of credibility framework & picture of credibility will be considered incomplete if not incorporated in the study.
		Political Astroturf Meme Detection	[16], [86], [87]	
		Spam and Phishing Detection	[88]–[91]	
		Topic specific Expert Identification	[92], [93]	
		Personality Specific Behavior Identification	[94]	
		Suspicious Behavior: Bot/Troll/Cyborg/Sybil/Content Polluter, Social Spambots, etc.	[80], [95]–[97] and [23], [26], [98]–[103]	
		Influence and Diffusion	[104]–[107]	
		Trust and Distrust Propagation	[108], [109]	
		Post Ranking	[110]–[113]	
		Hate Speech, Offensive and Abusive Language Detection	[114]–[116]	
		Hyper-partisan/Bias/Polarization Detection	[17], [102], [117], [118]	

186 **Information Credibility Taxonomy:** In the following sec- 193  
 187 tions, from section III to section VII, complete information 194  
 188 credibility is presented. The taxonomy is also drawn in figure 195  
 189 3. In this hierarchy the first branch named ‘Credibility and 196  
 190 its Components’ presents different types of credibility, cred- 197  
 191 ibility dimensions, , credibility constructs, credibility defini- 198  
 192 tions, etc. Second branch named ‘Theoretical Frameworks of 199

Credibility Assessment’, which actually presents evolution of Credibility, till date. In the field of communication and psychology such concepts are best presented as frameworks. In third branch named ‘Supported Research’ where different aspects of credibility i.e.: Deception, Hate Speech, and Influence Identification, etc are presented. Fourth branch named ‘Information Credibility of Social Media and Microblogs’,

S. No.	Category	Sub-Category	Reference	Description
4	Information Credibility of Social Media and Microblogs	User Perception	[37], [119]–[122]	Many organic surveys are conducted in which user perceptions or other elements have been studied, to explore all possible and important features of information credibility specifically with respect to the perception, judgement and heuristic of user.
		Explanatory Studies	[30], [120], [123]	Wide range of features are studied, and many explanatory studies are conducted regarding broad feature analysis. To conclude what serves best for credibility assessment data is collected from microblogging sites and tagged either by means of crowd sourcing environments or experts.
		Source Credibility	[19], [34], [124]–[130]	Researches where information credibility assessment is done through greater focus towards source /user of information
		Feature Based Models	[21], [131]–[139]	ML/IR based models are used which use features commonly related to Topic, Posts, Authors, and Network, etc. Either atomic level of information is used, means contents contained within the tweet or Varying level of information with aggregated and historic features, to assess the Information Credibility
		Graph Based Models	[128]	Uses SNA/ Graph based models by utilizing friends-followers network, user-tweet-retweet and retweet networks, etc.
		Hybrid Models	[31], [140], [141]	Some combination of Feature based and Graph/SNA based methods used
5	Standard Credibility Dataset	Credibility benchmarks are not predefined therefore its related gold standard dataset is missing. The difficulty of collecting large amount of such data has not yet received the attention it deserves [29].		

200 presents types of information credibility experiments, related 235  
 201 to social media and microblogs only. The last branch named 236  
 202 'Standard Credibility Dataset' presents details about avail- 237  
 203 able datasets. 238

### 204 III. CREDIBILITY AND ITS COMPONENTS 240

205 As an important objective with many challenges, this section 241  
 206 not only presents credibility definitions (as related terms) 242  
 207 but also extends them systematically and forms the basis 243  
 208 of the credibility framework's building blocks (e.g.: levels, 244  
 209 dimensions, constructs) through related research studies from 245  
 210 different fields. Different credibility components are compre- 246  
 211 hensively explored and presented. 247

#### 212 A. CREDIBILITY DEFINITIONS: 250

213 Many efforts have been made to define Credibility. It is a 251  
 214 complex and multi-dimensional concept. There is no clear 252  
 215 definition, it has been defined through several related con- 253  
 216 cepts [22]. Therefore such definitions are taken from both, 254  
 217 strong research studies and standard dictionaries:

218 It is defined as: "believability, trust, reliability, accuracy, 255  
 219 fairness, objectivity, and other concepts and combination" 256  
 220 [39], Oxford dictionary defines credibility as "the quality 257  
 221 of being trusted and believed in" [40], as Merriam Webster 258  
 222 dictionaries it is defined as "the quality of being believed" 259  
 223 [41]. Many researcher's core references of studies in commu- 260  
 224 nication examining credibility as message credibility, source 261  
 225 credibility, and media credibility [42]. The majority of re-  
 226 searchers are agreed that there are two attributes of credi-  
 227 bility: expertise and trustworthiness [22], [43]–[46]. Simi-  
 228 larly, across multiple definitions credibility is believability. 263  
 229 Credible information means believable information similarly 264  
 230 credible persons are believable persons [47]. 265

231 After going through the above formal definitions we can 266  
 232 divide credibility into two main components: message and 267  
 233 source. Where the source is further examined through trust- 268  
 234 worthiness and expertise. This forms the basis of credibility

framework.

**Credibility Components:** After an in-depth exploration of research studies conducted in psychology, communication and information science, and to understand the broad domain of credibility, the following major credibility-related components (e.g.: levels, dimensions, and constructs) are found. They all are comprehensively discussed in following subsections and summarized in table 2 and 3 as well. These components are in varying sizes/levels of hierarchy. The top most (levels of credibility) is defined first and the lowest most (constructs) is defined last. The order is also maintained in table columns. The outcome of the credibility components section would be resulted in section X and to some extent, section XI. The following components are explored from various studies to propose a generic credibility framework for social media. The framework simply exposes the relationships found in these components. In the last portion of section XI where generic social media framework is further transformed for microblogs, using microblog specific features is not concerned as an outcome of this section.

#### 249 B. LEVELS OF CREDIBILITY: 250

251 There are different levels of credibility assessed in literature, 252  
 253 which should be known for a better understanding of the 254  
 255 subject area. Levels of credibility are treated at the highest 256  
 257 level of the component's hierarchy or they are a macro- 258  
 259 level component. They are classified as following and also 260  
 261 summarized in table 2 and 3:

##### 262 1) Post Credibility: 263

264 It is the most important and primitive form. It means the mes- 265  
 266 sage or post itself is credible [136], [160]. It may effects the 267  
 268 credibility of the user or event, etc. It is the most suitable for 269  
 online/ real-time credibility identification systems because no 270  
 historic data is needed. On the dark side, it poses a weak 271  
 credibility assessment based on a limited scope.

**TABLE 2.** Credibility Components identified from research studies: different research studies related to Credibility Levels, Dimensions, and Constructs (table 1 of 2).

Ref	Levels	Dimensions	Constructs	Description
[142]	Source, Message	Quality, Trustworthiness	Source: Competence/ Expertise, Proximity/ Location, Popularity. Message: Recency, Corroboration/Agreement	Trustworthiness metrics proposed through survey research.
[143]	Topic, Source, Message	NA	Source: Authority/ Influence, Expertise, Popularity. Contents: Info. Quality, Popularity	Exploratory credibility feature analysis conducted on Twitter data, tagged by crowd-sourcing and experts
[144]	Source, Message	NA	Source: Expertise, Community. Message: Clarity, Emotions/Valance, Consensus (Consistency, User Judgment)	Social media based credible marketing related electronic word of mouth (eWOM) framework is proposed based on research theories.
[145]	Topic, Source, Message	Information Quality, Expertise, Trustworthiness	Survey covering many constructs used in studies.	Complete literature survey presenting different Levels, Dimensions, and Constructs of credibility.
[146]	Media Credibility	NA	1. Trustworthiness, 2. Un-Biased, 3. Accuracy, 4. Completeness, 5. Fairness	Defining and measuring media credibility.
[147]			6. Balanced (added)	Effects of balanced and imbalanced conflict story structure on perceived story bias and news media credibility explored through experimental study.
[148]			7. Factual (added)	Many constructs are measured through experimental study.
[149]			8. Expertise (added) 9. Social Concerns (added)	Literature review of credibility in the contemporary media environment.
[150]			Only 1-4	Survey on media credibility of newspapers accounts on Sina Weibo.
[151]	Source Credibility	Expertise, Trustworthiness	NA	Seminal work on source credibility: Survey & Controlled Group Study.
[152]		Goodwill/Caring (added)		First suggested perceived caring/goodwill as source credibility aspect.
[153]				Aspect of 'caring' fully studied in survey.
[154]				Reexamination of the construct and its measurement done and Goodwill added through survey study.
[155]				Endorsing through theories
[156]	General Credibility	Expertise, Trustworthiness	NA	Seminal work in Attitudes & Comm., reporting series of experiments on credibility.
[157]	Source, Contents	Quality, Expertise, Trustworthiness, Reliability/ Relevance/ Consistent	NA	Literature based, proposed contents/IR Credibility Framework

## 2) User/Source Credibility:

It corresponds to the poster (e.g.: speaker, organization, govt., news organization, etc.) or user of the post [126], [128]. In most studies, it is presumed that if the source is credible then the message associated with the user is also credible [34], [124], [130]. Somehow it is treated as the higher level, which means user credibility may be based on the user's post collections [34]. Which makes it a historic/ offline assessment system, because we need all historic data for evaluation. Online/ real-time or immediate assessment is not possible. Hence combined post and user information presents better credibility identification.

**Social/Domain Expert Credibility:** In [161] a variant or subset of source/user credibility is identified. It is based on the social status of a user in a social network on a certain domain. A similar concept is also used for Opinion Credibility [162]. Source credibility is known to be a super-set of such subsets.

Source credibility could be measured in terms of a broad set of credibility aspects like influence, popularity, truthfulness, expertise, biasness, etc. whereas such subsets are measured on just a single aspect e.g.: expertise.

## 3) Topic/Event Credibility:

Event comprises all related posts to a specific event/topic. Whereas topic/event could be identified by a set of keywords [31], [134], [163], [164]. The specific event comprises a collection of posts and associated posters as well. An example of such topic/event credibility is the Credibility of posts during COVID-19.

## 4) Media Credibility:

It is also multidimensional (high level) construct. Comprised of source credibility and medium credibility. Medium credibility focuses on the medium through which the message

**TABLE 3.** Credibility Components identified from research studies: different research studies related to Credibility Levels, Dimensions, and Constructs (table 2 of 2).

Ref	Levels	Dimensions	Constructs	Description
[67]	Credibility Constructs (Media, Source, Content)	NA	1. Believable/ Plausibility, 2. Truthful, 3. Trustworthy 4. Objectivity/Un-Biased 5. Reliability/ Accuracy/ Relevance/ Consistent	Unifying framework defined constructs
[68]			Found Best:(2-5 above) & 6. Recency/Timeliness, Found Good (for other Information Objects): 7. Completeness 8. Official, 9. Un-Biased, 10. Authority/Influence, 11. Expertise, 12. Scholarly/ Reference/ Educational Endorsement	Extension to Unifying framework to make it global
[158]	Content (Content Trustworthiness)	NA	1. Topic 2. Context and criticality 3. Popularity 4. Authority/Influence 5. Experience/ Reputation 6. Recommendation 7. Related Resources 8. Provenance/ Source 9. User expertise 10. Bias 11. Incentive 12. Limited resources 13. Agreement/ Corroboration 14. Specificity 15. Likelihood/ Believable/ Plausibility 16. Age/ Timeliness/ Validity 17. Appearance 18. Deception 19. Recency/Recent Image	Comprehensive study describing content trustworthiness: means how end-users make decisions regarding trusting information. Exhaustive literature review and simulation study supported.
[159]	Source, Message	Expertise: (Source, Content), Trustworthiness	Expertise: Quality, Accuracy, Authority, Competence Trustworthiness: Reputation, Reliability, Trust	Study from communication domain enlightening emergent and Modern concepts related to credibility.

is delivered (e.g.: newspaper, radio, television, etc.- In the context of our study it is just an underlying social network used for information propagation) [165]. In our case of microblog, the microblog’s credibility is Media Credibility which is based on the poster and underlying social network used for information propagation (as the medium). A very important and distinct notion presented in [59] that in modern scenario medium is also replaced with source only. Therefore only source (including all chain of message propagators) credibility could easily be used in place of media credibility.

The above types are somehow synchronized with each other. Therefore media credibility assessment system will require examination of the post, source, and underlying information propagation social network, to claim its microblog credibility system. Therefore for our proposed credibility framework only post-level and source-level credibility would be enough.

**C. CREDIBILITY DIMENSIONS AND CONSTRUCTS:**

It is quite challenging to define credibility in terms of its necessary components/elements, because there is no standardization due to its multidisciplinary [166] and emerging [159] nature. In the field of psychology and communication, the orientation of credibility is source-based and therefore called source credibility whereas in information science it is message oriented and called information credibility [166]. Dimensions are considered at middle and constructs are at the lowest level of credibility components hierarchy.

1) Credibility Dimensions:

Despite all above challenges it is observed through literature exploration that the majority of researchers accepts that there are at least two major dimensions (dimensions are also called topics, factors, etc. in literature) of credibility: Expertise

and Trustworthiness [151], [156], [159], other many studies endorse with minor addition [152]–[155]. Another important Dimension named: Information/ Data/ Content Quality is also found in [145], [157], [158], [167].

It could be concluded that the most agreed upon dimensions are Expertise, Trustworthiness, and Quality of Information. These could be the necessary dimensions of the proposed framework.

2) Credibility Constructs:

Under the above dimensions, there are some constructs (constructs are also called sub-topics, sub-factors, etc. in literature) proposed in different credibility studies. The list of constructs could be different concerning information object or media, etc. A very detailed survey discussing factors/sub-factors (topics/sub-topics) studied in variety of research studies [145]. Some basic credibility constructs are proposed in the most popular and highly concerned ‘unifying framework’ [67] (will be discussed in next section) which were extended concerning the varying type of information object (e.g.: Social Networks/ Media, Microblogs, Web Blogs, Search Engines, General Websites, Electronic Commerce Sites, News Sites, Educational Portals, etc.) or media contents (TV, radio, podcast, music, photo, video, etc.) in [68]. Detailed constructs specific to Data/ Content Quality are presented in [158], [167]. Constructs to assess Media Credibility are proposed in [146]–[149]

Regarding our proposed credibility framework which will be generic to social media but specific to microblogs. The levels and dimensions would be generic to social media only. Constructs must be compatible with both social media and microblogs and then further lower-level components (e.g.: features) must be microblogs specific or information object-specific only. Keeping the specific attributes of social media and microblogs both, the following few constructs



could be shortlisted from table 2 and 3 in addition to the following two criteria. 1. These constructs are common to both post and source levels, and 2. They are also common to trustworthiness, expertise, and information quality dimensions. These constructs are; 1. Recency, 2. Truthful, 3. Deception, 4. Topic, 5. Specificity, 6. Unbiased/Objectivity, 7. Popularity, 8. Plausibility, 9. Authority/Influence, 10. Competence/Reputation, 11. Uniqueness/Completeness, etc. Complementing the above recommended key Levels, Dimensions, and Constructs, some frameworks (comprised of levels, dimensions, and constructs) are developed and experimental studies are conducted to adhere to the findings discussed. For example, the electronic word of mouth (eWOM) framework for marketing related to social networks credibility is presented in [144]. The credibility framework for Information Retrieval systems is presented in [157].

In addition to the above frameworks and basic component related studies there are few exploratory studies conducted which also support and confirm the identified components. An exploratory study for credibility feature analysis conducted on Twitter data, tagged by crowd-sourcing and experts [143] (see table 2 and 3, for these frameworks) Summarized Levels, Dimensions and Constructs are presented in table 2 and 3. There are numerous studies found in psychology, communication and Information science on credibility-related components e.g.: levels, dimensions, and constructs; but only some representative studies are presented in the table for understanding and support.

#### D. RELATIONSHIP OF CREDIBILITY AND TRUST:

The concept of credibility and trust must be clarified and their relationship should be presented. Credibility and trust are mistakenly used interchangeably. Credibility is believability while Trust is dependability. Credibility is an antecedent to trust [48]–[51]

#### IV. THEORETICAL FRAMEWORK OF CREDIBILITY ASSESSMENT:

For the past many years, there have been so many research studies on credibility. All mostly in the field of information science, psychology, and communication. However, to better understand people's credibility assessment within various information contexts, modern credibility research has started to take a multidisciplinary approach [166] and becoming emergent [159]. In various research communities, different conceptual and theoretical frameworks have emerged regarding the conceptions of credibility, due to increasing concerns about the credibility of online information. There are the following distinct conceptual or theoretical frameworks categorized and described in order (similar to evolutionary generations), for examining the credibility of online information. They provide an understanding of the credibility assessment process and related concepts and how it is affected in general or discuss the underlying process involved behind people to

assess credibility. One can easily understand that how these frameworks are evolved concerning the modern requirements and challenges:

**4.1: Media-based Framework:** It is the earliest framework, developed within the field of communication. Researchers within this framework have long been interested, since the 1950s, to know the relative credibility [52] of different media channels (e.g.: Radio, TV, Magazine, Newspapers, and now Web is also included). Communication scholars investigated various factors affecting media credibility [53] including people's perception of Web-based information, and Web vs traditional media [54], [55].

The major limitation of this framework was that it considers people's general perception regarding medium instead of focusing on what use of information, which is obtained from it. For example, if someone considers the Web as the bad medium in terms of credibility doesn't mean that every website will be considered poor in credibility.

**4.2: Website-based Framework:** In this framework complete website is examined for credibility. In Stanford Web Credibility project [168] various elements of the website are examined which affects user's credibility assessments. After many studies Fogg's: Prominence Interpretation Theory is developed; which talks about the following, that needs to occur for people to assess web credibility: Prominence (likelihood of an element noticed) and Interpretation (value assigned to that element based on user's judgment). Factors affecting prominence as well as interpretation are also discussed [43]. There are few other studies [169], [170] found on website credibility under the website-based framework, all have the common strength that it covers both contents with peripheral cues (e.g.: appearance, design, presentation, etc.) as components of credibility. But on the other hand side, there is a weakness that every piece of information contained in the website is not separately considered.

**4.3: Content-based Framework:** Website contains many information objects therefore each information object is individually assessed in this framework. This framework assumes that information credibility may vary even within the same website. The main focus of the framework is: When we access any piece of information we emphasize assessing its quality. Therefore the chief aspect of information quality is defined as credibility [56]. It is reported in [57] that social-Q&A type of sites, users evaluate credibility primarily on contents because of having limited cues to source credibility. The weakness of the framework includes missing the emotional effects of interaction with information and aesthetic aspects of the information object.

**4.4: Interaction-based Framework:** This framework assumes that instead of discrete evaluative event credibility assessment is best expressed through an interactive and iterative process. It also guides that assessment of credibility

478 could easily be chalked out through observation during user's 534  
479 information seeking process with their selections made for 535  
480 searching that information. 536

481 The interaction framework also emphasizes the fact that 537  
482 credibility assessment is subjective means highly depends on 538  
483 the user's current knowledge and experience. Limitation to 539  
484 this framework seems that most of the studies only focus on 540  
485 the human information searching and navigating process. 541

486 542  
487 Rieh's model explains that when a user starts the information- 543  
488 seeking process, it begins earlier from predictive judgment, 544  
489 which leads the user to access information resources and 545  
490 then go towards evaluative judgment [58]. Hilligoss and Rieh 546  
491 added the third type of judgment as Verification [171], later 547  
492 through their empirical study. 548

493 Wathen and Burkell define an interactive and stage pro- 549  
494 cess where the first Website's surface-level characteristics 550  
495 (content organization, interactivity, interface design, speed, 551  
496 appearance, etc.)/medium credibility is rated, then the user 552  
497 rates the source and message (trustworthiness, competence, 553  
498 expertise, etc.) and the third aspect is the interaction of 554  
499 presentation and content [59] which is finally assessed as per 555  
500 user's cognitive states. 556

501 Sundar's credibility assessment also adheres interaction 557  
502 framework and presents the MAIN model (Modality, Agency, 558  
503 Interactivity, and Navigability) having four technical "af- 559  
504 fordances" in digital media [60]. Affordances can increase 560  
505 or decrease content effects on credibility, like moderators;  
506 in several psychological ways. It is therefore recommended 561  
507 by Sundar, that role of heuristics in credibility assessment 562  
508 should be explored. To understand the role of the heuristic in 563  
509 understanding credibility assessment is presented in Elabo- 564  
510 ration Likelihood Model (ELM) of persuasion and Heuristic 565  
511 Systematic Model (HSM) of information processing. Both 566  
512 models share many of the same concepts. Therefore Dual 567  
513 Processing model of information processing and credibility 568  
514 evaluation [66] has taken motivation into account like dual- 569  
515 process theories [172] and also based on both. 570

516 ELM of persuasion [61] is dual-process theory and the gen- 571  
517 eral theory of attitude change (e.g.: What attitudinal changes 572  
518 in user will occur when user come across messages and 573  
519 sources). It provides a general framework for understanding 574  
520 the basic processes underlying the effectiveness of persuasive 575  
521 communications. 576

522 Similarly, HSM of information processing [62] is a popular 577  
523 communication model which explains how people receive 578  
524 and process persuasive messages. Similar to all dual-process 579  
525 theories: ELM, Controlled and Automatic Processing Models 580  
526 (CAPM) [63], it is also defined in this model that individ- 581  
527 ual can process messages in either ways, systematically or 582  
528 heuristically. 583

529 Another widely used interpersonal communication and me- 584  
530 dia studies theory named Social Information Processing 585  
531 Theory (SIPT) [64], [65] which explains online interpersonal 586  
532 communication and how people develop and manage rela- 587  
533 tionships in a computer-mediated environment. It says that 588

the community exploits any piece of information that the  
channel provides them to make assessments about others.

Among dual-process theories (ELM, HSM, CAPM) and  
SIPT, there are few other fairly general theories and frame-  
works that are often adopted by credibility researchers to  
characterize the credibility assessment process and its con-  
structs and components.

**4.5: Unifying Framework:** Finally most important unifying  
framework of credibility assessment is proposed for a dif-  
ferent type of media, information objects, and contents for a  
variety of information activities. It provides very basic levels  
of credibility judgments: Interaction (credibility judgments in  
which sources or information examined), Heuristics (general  
rule of thumb, could be applied to a wide range of situations),  
Construct (how credibility conceptualized) as basic levels  
and an additionally defined Context (surrounding the user)  
of credibility assessment [67]. Later the framework was fully  
extended by Rieh et al. [68] to cater to the need of current  
and modern participatory web environment (include Web  
2.0 means all kinds of modern social media services and  
others). It could be concluded that Unifying Framework  
is the most relevant and therefore should be followed to  
fulfill the modern requirements. The proposed framework  
is also enriched with the constructs presented in Unifying  
Framework.

## V. SUPPORTED RESEARCH:

Many of the supported or closely related and somehow  
different dimensions of microblogs-based information cred-  
ibility, have already been studied separately. Unfortunately,  
they are not considered as directly related to credibility  
in the literature, but all of them are comprising different  
constructs/aspects of credibility and therefore need to be aug-  
mented, holistically. The mapping of all supported research  
studies with appropriate constructs is done in this section.  
All these constructs/aspects are also shown in the proposed  
high-level credibility framework's table:14 and then these  
aspects are mapped to individual features in table:15, where  
all these studies are highly contributing. We consider these  
supported research studies as important building blocks of  
credibility or necessary components of the credibility frame-  
work. Picture of credibility will be considered incomplete  
if they are not incorporated in the study. Each one of them  
is considered a completely separate research area therefore  
details are omitted but only the research area name together  
with important references are mentioned. Important terms are  
defined in the Appendix for basic understanding and clarity.  
What we have done for simplicity and increased productivity  
that we go through all supported research studies and list  
down all important features. These features are then proposed  
for implementing microblogs specific credibility framework.  
They are presented in the table: 15 which provides the im-  
plementation of our generic framework to microblog specif  
framework. All these features are added with their supported

589 references and reason in table 15 of our proposed credibility 645  
590 framework section XI. 646

591 In this section, to support the understanding of credibility 647  
592 components, each area of research (which is named as 648  
593 supported research in the study) is categorized with respect 649  
594 to its respective level and appropriate construct. For exam- 650  
595 ple, 'fake news detection' is an area of research which is 651  
596 classified under construct, named 'Deception and Truthful', 652  
597 presented within the level, called 'post level'. The area of 653  
598 research will not be discussed, only the name of the area 654  
599 with respective references will be included. There are few 655  
600 constructs shortlisted in section III-C2 related to social media 656  
601 and microblogs-based information credibility. Examples of 657  
602 those few constructs are; 1. Recency, 2. Truthful, 3. De- 658  
603 ception, 4. Topic, 5. Specificity, 6. Unbiased/Objectivity, 659  
604 7. Popularity, 8. Plausibility, 9. Authority/Influence, 10. 660  
605 Competence/Reputation, 11. Uniqueness/Completeness, etc. 661  
606 Different areas of research considered related to this section 662  
607 are categorized under these relevant constructs. Those areas 663  
608 of research under each construct's heading are as follows. 664  
609 The constructs are also grouped under respective levels like 665  
610 post level and user level. 666  
611 667

612 **Post Level Constructs:** It is discussed earlier that post is the 668  
613 most basic and lowest level in all other levels of credibility. 669  
614 Though we have considered only two levels, post, and user. 670  
615 Aggregation of many post-level constructs will automatically 671  
616 result in user-level constructs, e.g.: if the majority of posts 672  
617 are biased then the user will automatically be biased. The 673  
618 same will be the case of fake posts. It means that few 674  
619 constructs will be common in both levels. Those common 675  
620 constructs are Deception, Truthful, Unbiased, and Popularity, 676  
621 etc. Despite that few constructs are common, only those 677  
622 constructs are repeated where detection mechanism found 678  
623 different at both post and user levels, e.g.: Deception and 679  
624 Truthful. The techniques detecting deception at post level are 680  
625 discussed as fake news detection, rumor detection, etc. but 681  
626 techniques detecting deception at the user level are called 682  
627 bot-detection, suspicious behavior detection, etc. 683  
628 684

629 **5.1: Deception and Truthful:** Detection of all deceptive 685  
630 and untruthful contents must be done at each post level. 686  
631 This section includes all such studies which provide the 687  
632 understanding and also suggest ways and means of their 688  
633 detection. 689

634 It is discussed in [173]–[175] that false Information [83], 690  
635 [176] or deceptive information [33] has variety of flavors: 691  
636 Fake/ False News, Misinformation, Disinformation, Hoaxes, 692  
637 Propaganda, Satire, Rumors, Click-Bait, and Junk News, etc. 693  
638 Though an agreed and standardized definition is completely 694  
639 missing but is generally considered that misinformation is 695  
640 information that is inaccurate and misleading which could 696  
641 spread unintentionally in contrast to disinformation which is 697  
642 false information and spread deliberately to deceive people. 698  
643 699

644 False Information Detection: Following are studies related 700

to deception and false information detection including their  
different forms. Only name of the field/area and related  
references will be provided.

Misinformation/ Disinformation and its detection: [20],  
[174], [177]–[180], Rumor and its detection: [3], [69]–  
[76], Fake News and its detection: [25], [77]–[85], Stance  
Detection is basically identification of the relevance of news  
article's contents with title. Its now assumed as sub-category  
of fake news detection, such that for fake news identifica-  
tion first stance is evaluated: [181], Hoax Detection: [182],  
Spam and Phishing Detection: [88]–[91] spam and phishing  
detection techniques can also automatically filter click-bait,  
fake reviews, and some political astroturfs. Because they are  
similar in structural or strategical patterns and may called  
modern form of spams.

Damage of Reputation Detection: There are some types of  
deceptive and false information that damage one's reputa-  
tion and naturally affect one's credibility, they are called  
smear campaigns which may include: satire, conspiracy,  
propaganda [183], political astroturf memes, etc. There are  
different Political Astroturf Meme Detection studies also  
found: [16], [86], [87].

**5.2: Bias/Objectivity:** It is found that some post may have  
a piece of such information which come from a particular  
point of view and may rely on propaganda, decontextualized  
information, and opinions distorted as facts. These posts are  
categorized as extremely biased. They must be identified  
or detected in the early stages of their spread otherwise  
have associated grave repercussions. They create highly  
polarized groups, in terms of religion, politics, race, etc.  
Therefore following are few example studies which can  
identify 'bias/objectivity' construct of credibility, they are:  
Hyper-partisan/ Bias/ Polarization Detection: [17], [102],  
[117], [118].

**5.3: Plausibility/Likelihood:** Freedom of expression is a  
human right but hate speech towards a person or group  
based on race, caste, religion, ethnic or national origin,  
sex, disability, gender identity, etc. is an abuse of this  
sovereignty. Hate speech is essentially a discourse that might  
be extremely harmful to the feelings of a person or group  
and may contribute towards brutality or insensitivity which  
shows irrational and inhuman behavior. It seriously promotes  
violence or hate crimes and creates an imbalance in society  
by damaging peace, emotions, reputation, trust, credibility,  
human rights, justice, and democracy, etc. In addition, to  
hate speech some other related concepts must also be con-  
sidered like Hate, Cyberbullying, Discrimination, Flaming,  
Harassment, Abusive Language, Profanity, Toxic Language  
or comment, Extremism, Radicalization, etc [14]. These all  
are some general information quality-related constructs that  
must be considered for detection. Following are few example  
studies which can fulfill the requirements, such as; Hate  
Speech, Offensive and Abusive Language Detection: [114]–

- [116].
- User Level Constructs:** User level is higher than post level. Many user-level constructs could be accumulated through their respective post-level constructs. Therefore they are omitted from this section. Considering the case of fake posts, if the majority of posts posted by a user are fake then that user will not be trustworthy. In the following user-level constructs, only those constructs are presented where detection mechanism is found different concerning the user. Following are all supported research studies categorized under user-level constructs. Only the name of the field/area and related references will be provided, details of the field are not included:
- 5.4: Deception and Truthful:** It is worth mentioning that majority of the incredible contents are spread through different types of Bots, Trolls, Cyborgs, Sybils, Content Polluters, or Social Spambots, etc. There are almost 15-17% accounts which are bots [23], presenting human impersonation and perform many malicious and suspicious activities, e.g.: Spread of misinformation and fake news, fake support, fake product reviews, advertise for doubtful legality, hashtag and other promotions, spread unsolicited spam, scam URLs, terrorist propaganda, manipulate the stock market, rumor dissemination and support, conspiracy, astroturf political campaigns, and religious activism, bias public opinion, sponsor public character and many similar activities [20], [26], [95]. Therefore once they are identified and blocked then all such contents will automatically be filtered and the remaining large portion of contents will be treated as legitimate and credible.
- There are studies found for such malicious profiles identification and detection, for example Bot/ Trolls/ Cyborg/ Sybils/ Content Polluters/ Social Spambots and its detection: [23], [26], [98]–[103] and Suspicious Behavior Detection: [80], [95]–[97].
- 5.5: Competence, Topic and Specificity:** Following are a few examples of supported research studies that could help in the determination of the above group of constructs. Dealing with the social status of a user in microblog’s social network on the certain domain such as politics, education, sports, science and technology, social issues, etc. It is simply called Topic Specific Expert identification, here the competence within a specific domain or topic is concerned, which could be done with the help of these studies: [92], [93], [161]. A very similar concept is used as Opinion Credibility in [162]
- 5.6: Popularity, Influence/Authority:** Every user in a microblog’s social network has certain influence/ authority/ popularity. Highly influential/ authoritative/ popular users can affect an individual’s attitudes, beliefs, and subsequent actions or behaviors. We need to identify an appropriate way to measure user influence/ popularity/ authority score. Most authoritative/ influential/ popular users are assumed
- more credible. There could be different ways of measuring such scores. Making use of the follower-following network or user-tweet/retweet network, and then apply modified page rank like model or some form of authority transfer or some centrality measure for calculating highly influential/ popular/ authoritative/ reputed user. It could be measured by applying some ratios of followers count, followings count, with some form of popularity measures e.g.: no of times a user is mentioned, retweeted, replied, listed, favorited, etc. by other users of microblog’s social network.
- The above methods are commonly considered in computing source credibility. Some good variants can enrich these methods with a quite different perspective. Using the following concepts will provide required value addition, such as Post Ranking, which is done concerning relevance of user and content, as well as source popularity: [110]–[113]. Influence and Diffusion Methods: [104]–[107]. Trust and Distrust Propagation: [35], [108], [109]. Personality specific behavior [94] identification, which greatly helps in detecting different behaviors. Personality Detection, which provides big-five personality traits (i.e.: 1. Open/Closed, 2. Spontaneous/ Conscientious, 3. Introvert/ Extrovert, 4. Hostile/ Agreeable, 5. Stable/ Neurotic) that help predicting behavior and influencing ability. [184], [185].

## VI. INFORMATION CREDIBILITY OF SOCIAL MEDIA & MICROBLOGS

The outcome of our study is two-fold. Understanding the broader domain of credibility with basic components identification and then the development of compatible social media generic framework. This will further be transformed to microblog-specific implementation. Considering the first objective: credibility related various generic studies from different fields have already been explored in former sections (III-B, III-C and IV). Presenting frameworks, models/theories (see section IV), and its macro components ( see sections III-B, III-C, e.g.: levels, dimensions and general constructs) for broad range of information objects (e.g.: General Websites, News Media Sites, Search Engines, etc.). Moving forward towards the second objective: it is needed to exhaust only information object-specific studies. Characteristics of our information object (social media and microblogs) are quite distinct from other information objects, such as the authenticity of the source is hidden from the user. Contents are massively shared. User engagements and responses are shown. Content has a long propagation path that is hidden from the user. User-generated content, which is noisy. Having spelling mistakes, free from grammar, small in size, have little context, contain language variations, furnished with special meaning in form of emoticons, hashtags, user mentions, re-tweets, and capitalization, etc. Therefore we have only considered social media and microblogs specific studies in this section. Considering any other type of information object (e.g: General Websites, News Media Sites, Search Engines, etc.) related credibility studies will not be productive for this section. Credibility constructs are somehow information

object-related and need transformation [68] which is done in many studies, like [186]. It is also done in our proposed credibility framework presented in section XI, where constructs are only social media-specific and then corresponding features are microblogs specific.

There are some domain-specific studies of information credibility found, like: Health [29], Disaster [187], Fake Review/Opinion [188], Image/ Media [189], Geographic Information [190], Language Specific [191], [192], Country-Specific Perceptions [122], etc. All such studies are not considered much relevant because the challenge here is to understand correctly what is information credibility concerning social media in general first and then how will it be achieved for microblogs. Once this general understanding of information credibility will be developed then these very specific studies will be effective.

The outcome of this section has resulted in the last segment of section XI as well, where the generic credibility framework of social media is further transformed to microblogs. Studies of this section guide that what are the set of microblog's features which are recommended for specific aspect (e.g.: Hate, Bot, Fake, Influence, etc.) evaluation.

Studies conducted specifically on information credibility related to social media and microblogs can easily be classified as studies that present User Perceptions of Information Credibility, Explanatory Studies, Source Credibility, Feature-Based Models, Graph-Based Models, and Hybrid Models of Information Credibility.

#### A. USER PERCEPTION:

This section presents extremely important, Social Media and Microblogs Credibility specific variety of hypothesis related to human cognitive heuristics, judgments, perceptions, and assessments. Which are identified and examined through different methods like surveys, interviews, empirical & experimental studies, observations, and statistical methods. All such studies are very comprehensively presented under different columns in table 4 and 5. In each study of the table, a very organic survey is conducted in which user perceptions or other elements have been studied, to explore possible and important features of information credibility for social media and microblogs, specifically concerning the perception, judgment, assessment, and heuristic of the user. Researchers use these recommended features as a starting point and conduct explanatory studies to conclude what serves best for credibility assessment.

#### B. EXPLANATORY CREDIBILITY STUDIES:

There is no accepted credibility standard [20], [21] and it is very difficult to judge different researches and generalize the findings. In this section, such studies are included in which different efforts have been made for identifying important microblogs specific credibility indicators through the wide range of factors studied (see studies conducted in section VI-A to explore important credibility indicators), and then these explanatory studies are conducted. They conclude that

what serves best for credibility assessment. In such studies, to provide detailed features exploration and analysis, mostly data is collected from microblogs sites and tagged either through crowd sourcing environments or experts. A complete list of explanatory studies is presented and summarized under different columns in table 6. Following are only a few studies discussed for a basic understanding of such studies.

In [30] manually tagged dataset having three classes of features: social, context, and behavioral are analyzed, within 8 different topics and concluded the best credibility indicators. In [123] an effort has been made and a wide range of factors are studied and an explanatory study is conducted.

Another very important explanatory study together with user perceptions has been conducted in [120], which examines the relationship between reader's demographics and related credibility features with user perceptions. Over 1317 attributed news tweets were collected and annotated using both TweetCred and manually; for examination of the relationship between eight tweet level features (including source) having reader's perception of credibility, news attributes, and reader demographics features. Further correlation among the attributes was also explored using Cohen's Kappa, chi-square, and association rule mining.

#### Microblogs based Automatic Credibility Assessment Models:

Following are all such categories in which only microblogs-based automatic credibility assessment systems are considered. They are classified as; Source Credibility, Feature-Based Models, Graph-Based Models, and Hybrid Models of Credibility. The outcome of this section is resulted in section IX and section XII. In section IX all these automatic assessment studies are summarized in four groups and their important findings are discussed. Findings include common features, strengths, and shortcomings. In section XII recommendations are presented, based on important findings of section IX.

#### C. SOURCE CREDIBILITY:

There are many research studies where information credibility assessment is done through greater focus towards source /user of information [34], [130]. Ranking microblog users regarding their credibility could also be a candidate approach [124], therefore ways for source determination is also studied [125] and what affects the source credibility [129]. For example, in [127] US Senate voting history data is used and the user is ranked to measure information credibility based on their online behavior. CredRank Algo (based on IR tech.) is developed by the authors to detect Coordinated Behavior. If it is found then those users were marked as not-credible. In [19] researchers proposed that user influence can be measured through characteristics like In-degree, Retweets, and Mentions.

Focusing on source credibility, tweet timelines of 10 general and 10 highly influential Twitter users of five areas each like: car, investment; are fetched and then making use of

**TABLE 4.** Following are many surveys conducted. In these surveys user perceptions, judgment, heuristic, assessment or other elements have been studied. These elements are identified and examined through different methods like: surveys, interviews, empirical & experimental studies, observations and statistical methods (table 1 of 2).

Paper	Features Level Covered	Approach	Technique	Variables/ Features	Remarks
[193]	Topic, Post	Survey: Amazon MTurk	Variance Inflation Factor(VIF), Correlation, Hierarchical Regressions, Cronbach's $\alpha$	NA	Social media sites vs traditional news media
[194]	Topic	Survey: Online	Cronbach's $\alpha$ , 9 Hierarchical Regressions	NA	Politically interested online users view for social networks as credible
[195]	Post	Survey: Amazon MTurk, Fluo, Apollo	Maximum Likelihood Estimation, ANOVA	NA	Human cognitive limit vs effect of automated system recommendation
[196]	Topic, Post, User	Survey: Online	Variance Inflation Factor (VIF), Correlation, Hierarchical Regressions, Cronbach's $\alpha$	NA	Political blog credibility and selective exposure, avoidance
[121]	User, Post	Survey: Online	Statistical Methods, ANOVA, MANOVA	NA	Effect of follower-following over source credibility
[197]	Topic, Post, User	ACT-R Model Memories	Correlation, LDA, ANOVA	NA	Human credibility judgments
[8]	Post	Interviews, categorization using content analysis	Empirical Study	NA	Audience aware credibility constructs
[198]	User, Post	Survey: Amazon MTurk	Correlation Analysis, Statistical Methods	22	Factors influencing credibility perceptions for micro-blogs.
[199]	User, Post	Credibility Judgments	Cognitive Heuristics	NA	Cognitive heuristics for credibility judgment in online environments
[200]	Topic, Post, User	Survey: Mock Site, Interviews, Three Experiments	3 Way ANOVA, Cronbach's $\alpha$	5	Factor influencing credibility of health and safety information on Weibo
[201]	User, Post	Controlled Experiment of 2 Treatment Group	1 Way K-Group MANOVA, ANOVA	7	Twitter's human agent vs bots
[122]	Topic, Post, User	Survey: Online	ANOVA	5	Country specific credibility perceptions
[202]	User, Post	Empirical Study, Web-based information activity diary survey, Experience Sampling	Statistical Methods	11	Various credibility constructs
[203]	User, Post	3 Surveys using Mock Site	Post Hoc Wilcoxon Rank, Omnibus F, Tukey, Friedman's Test, 1 Way ANOVA, PCA	6	Social network derived credibility

921 topic-related user's social structure, they try to find most 935  
922 influential/centric users within each topic as credible [126]. 936  
923 It is the combination of topic models over message contents 937  
924 and link structure analysis of the underlying social network. 938  
925 User ranking based on authoritative user scores considering 939  
926 friend network and user-tweet/retweet network is imple- 940  
927 mented using ObjectRank in [128]. 941  
928 The research study performed in [112] focused on exploring 942  
929 indicators of credibility during eight diverse events. They 943  
930 concluded that URLs, Tweet length, Mentions, and Retweets 944  
931 are the best credibility indicators. The system proposed a 945  
932 ranking strategy based on content relevance and account 946  
933 authority considering: followers, mentions, list membership, 947  
934 and user-retweet graph. The system was trained using a 948

learning to rank algorithm named RankSVM.

In short: the majority of researches in this category make use of the follower-following network or user-tweet/retweet network, etc.; with some form of popularity measures e.g.: the number of times a user is mentioned, retweeted, replied, listed, favorited, etc. by other users of the social network, and then apply modified page rank like model or some form of authority transfer for calculating highly influential/popular/authoritative/reputed user as credible. In source credibility identification, Social Network Analysis (SNA)/Graph-based methods are exploited most of the time, except few studies, which found some weighted ratios of a different combination of popularity measures, effective. There is no consideration towards post quality therefore

**TABLE 5.** Following are many surveys conducted. In these surveys user perceptions, judgment, heuristic, assessment or other elements have been studied. These elements are identified and examined through different methods like: surveys, interviews, empirical & experimental studies, observations and statistical methods (table 2 of 2).

Paper	Features Level Covered	Approach	Technique	Variables/ Features	Remarks
[204]	Post	Survey Embedded Experiment	Statistical Methods, Regression, Cronbach's $\alpha$	11	Credibility of news: source, context
[205]	User, Post	Survey: Online	Kruskal-Wallis, Mann-Whitney U, the Wilcoxon Matched Pairs Test, Spearman Rank Correlation, Cronbach's $\alpha$ , Statistical Methods	NA	Credibility perception of social, teacher, scholarly tweets
[150]	User	Survey: Online	Correlation, P-Value, T-Test	4	Media credibility of newspapers accounts on Sina Weibo
[206]	User	Data Collection: Twitter, Coded: Research Team	Statistical Methods	4	Tweet's source credibility : fukushima nuclear disaster
[207]	Topic, Post, User	Tagging: Professionals, Features Rated (Perception Based): Amazon MTurk Survey	Krippendorff's $\alpha$ , Pearson correlation, Box Plots, Scatter Plots, P-Value, Precision, Recall, F1 Scores	6	Epistemic study of information verification: features for Hurricane Sandy pictures real/fake
[119]	Topic, Post, User	Think aloud, Elaborative Questions (Verbal)	Statistical Methods, ANOVA	31	Microblog credibility perceptions
[208]	Post, User	Survey: Online	Tukey's HSD Test, Hierarchical Regression Model, Constant-Comparative Method, Statistical Methods	3	Student perceptions of instructor credibility and beliefs about Twitter as a communication tool
[209]	Post, User	Two Software based Surveys	Linear Regression, Statistical Methods, Pearson Product Moment Correlation	5+3	Visualization perception of five + three factors of trustworthiness
[210]	Topic, Post, User	Survey Questions/Ratings: Office Users	Pearson Correlation, KMeans, Linear Regression, Feature Distributions, T-Test, Density Estimation: Gaussian Kernel, Outlier: K-Divergence, Statistical Methods	10	Study bias amongst microblog users due to the value of an author's name.
[211]	Post	Survey	Maximum Likelihood Estimation, Structural Equation Modeling, Statistical Methods, Error Methods: Chi-Square, RMSs, GFI, CFI, AGFI, CI.	8	Credibility and trust in online media use
[212]	Post, User	Survey: Online	Tool: G-Power, Paired Sample T-Test, ANCOVA, Wilks' Lambda, Levene's Test	6	Journalistic credibility on twitter

949 labeling the post as credible or not credible is completely 960  
 950 ignored. Primarily the efforts are being made to rank the user 961  
 951 therefore there is strong overlap with both ranking and graph 962  
 952 base credibility assessment methods.

#### 953 **D. FEATURE BASED MODELS FOR CREDIBILITY:**

954 Studies in this category usually build models which are either 966  
 955 Machine Learning (ML) based or Information Retrieval (IR) 967  
 956 based. They use features related to 'topics', 'posts', 'authors', 968  
 957 'network', etc., and of different types as well, such as 969  
 958 Aggregated and Historic. Examples of topic-level aggregated 970  
 959 features are the number of positive sentiment tweets, Avg. 971

length of a tweet in a topic, etc. Historic features are difficult  
 to extract and next level to aggregated features. For example,  
 A user will be known as Topic Expert if his number of tweets  
 under that topic is greater than the average number of tweets  
 of that topic tweeted by all users. Calculating such features  
 requires exhausting the complete dataset for that feature  
 level (e.g.: user in this example). Such types of features are  
 used to explicitly exploit inter-entity relationships, which are  
 inherent in graph/ network.

These feature-based assessment studies are also summarized  
 in tables: 7 and 8. Each study is comprehensively presented  
 across many important attributes. They are salient qualitative

**TABLE 6.** Explanatory Studies: Many efforts have been made for identifying important microblogs specific credibility indicators through wide range of factors studied in previous survey studies.

Paper	Features Level	Approach	Technique	Variables/Features	Remarks
[213]	Topic, Post, User	Tagging: CrowdFlower	Predictive Association Rule Analysis	8	News related tweet's credibility perception
[30]	Topic, Post, User	Tagging: Amazon MTurk	Distribution Analysis	34	Credibility related features distribution of twitter
[123]	Topic, Post, User	Tagging: Amazon MTurk	Statistical Methods, Kappa-statistic, Correlation, Forward Subset Selection Regression (FSS), Logistic Regression	45	Twitter credibility feature exploration and various ground truth analysis
[214]	Topic, Post, User	Tagging: Amazon MTurk	Statistical Methods, Kappa-statistic, Correlation, Forward Subset Selection Regression (FSS), Logistic Regression	45	Twitter feature exploration with network context and ground truth selection for credibility
[215]	Topic, Post, User	Tagging: CrowdFlower, Author and Post	Mean, Pearson Correlation	5	Impact of author's location on credibility
[216]	Post, User	10M Tweets Rated using proposed equations	Correlations, CDF, Statistical Methods	18	Scored features are statistically explored for trustworthiness assessment
[3]	Topic, Post, User	Manually (keyword search) Tagged for Ground Truth	Descriptive statistics, Filter Based Heuristic Approach	6	Understanding rumor/fake patterns/behavior/features in crisis
[143]	Topic, Post, User	Credibility Rating: Crowdsourcing and Experts	Krippendorff's $\alpha$ , Feature Distributions, Statistical Methods	44	Determining features of credibility in Arabic microblogs determining credibility
[120]	Topic, Post, User	TweetCred: Rating, CrowdFlower: Perception Survey	Chi-square Correlation Analysis, Cohen's Kappa, Association Rule Mining	Twitter:11, Demographic:4	Perception of reader vs news related microblog credibility features

972 dimensions of these research studies which should be known  
973 for efficient exploration of the research area. These attributes  
974 are as following:  
975 1. Paper: provides the reference of the concerned study.  
976 2. Algorithm: provides the name of the best performing algorithm  
977 of the study. 3. Learning Type: presents what type of learning  
978 or method is used like supervised classification, semi-supervised,  
979 supervised, unsupervised, ranking, etc. 4. Approach: presents  
980 what type of approach is used like feature-based, graph-based,  
981 information retrieval (IR) similarity measure based, weighted  
982 equations for scoring, user-defined ratios, etc. 5. Features Level:  
983 specifies that what different types and levels of features are used.  
984 There could be different levels of features e.g.: topic, user, post/  
985 tweet, or if its graph-based method then what type (directed,  
986 undirected) of the graph is developed over what entities/nodes  
987 (topic, user, post). Similarly, there are different types of features  
988 like historic, aggregated or temporal. User+Historic means that  
989 user-level historic features are used. 6. Dataset: shows summary/  
990 statistics of data collected in the study. All of the studies extract  
991 their own dataset, because of the unavailability of the standard  
992 dataset. 7. Outcome: what was predicted in the study is expressed  
993 in the outcome. e.g.: credible (credible, not-credible), credibility  
994 levels (high, medium, low, not-credible), rank/score (0-10), etc.  
995 8. Label Method: provides that who labeled the data like domain  
996 experts, crowd source workers, automatically tagged through  
997 computations, by authors, evaluators (manual team working for  
998 data extraction and labeling), manual (means labeling source is  
999 not defined). The labeling method defines the quality of the  
1000 system. The best labeling is done by experts while labeling done  
1001 by crowdsource workers is weak.

9. Focus: exposes major and special focus of the study, or system  
tile, e.g.: real-time assessment system, if the system is developed  
for 'emergency situation', if the system is produced for 'high  
impact events', 'fact-checking and scoring system', 'topic  
credibility', etc. 10. Product: either study is providing product  
as 'browser plug-in' or 'Twitter plug-in', or its just a research.  
11. Distinct Attribute: it provides highlights of the study or  
some distinct features of the study, or if the system uses some  
distinct components, or some methodology, like 'online  
emergency monitoring component' is provided, 'Experimental  
study' is also provided, post 're-tweet network' is exploited  
for assessment, 'topic-based' method is provided, the system  
explicitly works on 'user expertise and reputation' for  
assessment, the system provides idea of how 'topic-based  
expert user with biasness' is assessed, etc. 12. Category:  
this attribute provides fine-grained classification of the  
study, either system uses ML, or IR, Learn to Rank  
(ranking), Mathematical, or Hybrid methods for  
assessment.

There are generally two classified groups where first  
includes such studies in which scientist worked at the  
atomic level of information means only or mostly on  
tweet [136], [160]; to assess the Information  
Credibility, such as: In [160] it is assumed that  
credibility can be judged from tweet text, a credible  
tweet always has many retweets with original text  
remain. However in a low credible message several  
terms are added with user opinions, deleted or  
edited, and has low retweets. Based on the said  
concept user credibility is



**TABLE 7.** Feature Based & Hybrid Microblog Credibility Assessment Models Classification ( table 1 of 2). Each study is comprehensively presented across following important attributes. They are salient qualitative dimensions of these research studies which should be known for efficient exploration of the research area.

Paper	Algo	Learning Type	Approach	Features Level	Dataset	Outcome	Label Method	Focus	Product	Distinct Attributes	Category
[138]	Feature Rank Naïve Bayes	Classification	Feature Based	Tweet, Aggregated (Topic, User), User+ Historic	Yemen Civil War: Keywords (Taiz & Aden) Tagged 11000 sample tweets	Credible	Experts	Four Component System, CDF Based Feature Analysis	No	User Expertise & Reputation	ML
[135]	Proposed CIT Bayesian Network	Classification	Feature Based	Tweet, User + Aggregated, Topic+ Aggregated, Temporal	UK-Riots related topics, 350 Tweets Tagged	Credible	Experts	Emergency Situation	No	Online Emergency Monitoring Component	ML
[136]	SVM Rank	Semi-Supervised Ranking	Feature Based	Mostly Tweet Level	6 High Impact Crisis Events (tagged 500 tweets for each event)	Score/Rank	Crowd Worker	Real-time Credibility Score	Browser Plugin	Mostly features extracted from tweets	Learn to Rank
[134]	J-48 Decision Tree (3-Fold)	Classification	Feature Based	Tweet, User, Topic + Aggregated, Retweet Tree Propagation	2524 Trending Topics, Classes: News +Chats, Topic are tagged using 10 tweet samples	Credible	Crowd Worker	Topic Credibility	No	Retweet Net. tree used	ML
[132]	SVM Rank with PRF Re-ranking	Ranking	Feature Based	Tweet, User	14 News events of finance, politics like domains, 3586 Trending topics, 35M tweets + 6M Users, Tagged 500 Tweets for each topic	Rank	Crowd Worker	High Impact Event	No	Pseudo Relevance Feedback for Re-ranking	Learn to Rank
[133]	Similarity Score TF-IDF	Un-Supervised	IR Similarity & Feature Based	Tweet, User	2 News Topics (Iran/Yemen & Houthi), 600 Arabic Tweets, 179 Authorized News Articles for Verification, 29 Tagged Tweets for Evaluation	Credibility Levels	29 Tweets Tagged for Evaluation & Thresholds are set by Experts to Classify	Similarity Based News Fact checking Score+ Feature Score: Link, User Authority, Inappropriate Words	No	Experimental Study	IR Based
[139]	Random Forest (n-estimator:100)	Classification	Feature Based	Tweet, Topics, User + Historic	10 Trending Topics of Japan, 200 Tweets of each Topic Tagged	Credible	Crowd Worker	LDA used for Topic Extraction	No	Topic Based User Expertise & Bias Extracted	ML
[131]	J48 Decision Tree	Classification & Weighted Equations	Feature Based, Social Context Based Authoritative User using Ratios	Tweet, Topic+ Aggregated, User+ Historic	7 Topics of Libya only, 37K Users, 126K Tweets, Tagged 5000 only	Credibility: +Ve, -Ve, True, Null	Evaluators	Topic Based Assessments, LDA is used for Topic	No	1. Social Model: Social NW's important indicator's ratios are weighted. 2. Supervised Model based on Contents. 3. Hybrid 1& 2.	ML & Mathematical

also calculated with tweets. The reputation-based credibility degree assessment method developed for wikis is applied to tweets. The study has no experiments and Evaluation. It just uses ratios/ mathematical scores.

A browser plug-in named TweetCred, is a real-time system build over semi-supervised learning using SVM-Rank and trained through 45 tweet level features (only data provided in tweet object is used). These features are generally classified in Meta-data, Content-based linguistic features, Author (#follower-following and age), Content-based lexical features, URL reputation score, and Tweet Network features. The system is developed through six high-impact crisis events of the year 2013. Only US-based annotators were used to annotate 500 tweets. The system was widely downloaded and used [136].

Other group includes studies which exploit other level features too in addition to tweet level with all features types e.g.: Aggregated and Historic, such studies include A Hybrid model combining two models through averaging and filtering: the first model, named social model measure social credibility, deals with credibility at the user level, combining many dynamics of topic-specific content flow within its social network; and second model named content model measure content credibility, calculates fine-grained tweet level content based credibility [131]. In short: a total of 19 features are used to generate a score first and then making use of user friendship network user transfer that score to their followers. Dataset was generated through 7 topic-specific "Libya" and a total of 5000 manually annotated tweets of 37K users.

14 high impact news events of 2011 are considered and investigate the tweets based on supervised learning, with RankSVM + Pseudo Relevance Feedback over content-based and user-based static features, and then credibility is ranked [132].

An experimental system was developed with two approaches. One was based on the similarity of tweet news text and verified/authentic news text, and the other was combined with similarity-based features and other proposed (tweet and user level) features. Only IR-based methods were used and the system was developed on two hot news topics having 600 tweets which were verified through 179 authentic news articles [133].

It's a seminal study [134] where tweets belong to trending topics are collected and a wide variety of features related to Topic, User, Propagation, and Message; are extracted for supervised learning using J48 Decision Tree as best ML Algo. [135] Aims to measure credibility in an emergency situation using Bayesian Network over features based on: Diffusion, Topic, Content, and User.

An effort was made to develop a time-efficient twitter plugin in [137]. A dataset of 7000 tweets fetched on Nature Environment Preservation, with the help of more than 100 related terms, then 1206 tweets tagged and Random Forest classifier was trained over user and tweet level features. Results were improved through a reconciliation system for

tagging evaluation and re-tagging.

A classification system consisting of four components: Reputation Component - based on user popularity and sentimentality; it initially helps to filter neglected information for further assessment. Classifier Component - classify credible/incredible, using four ML-based classifiers. User Expertise Component - rate user expertness for the topic. Finally, the Feature Rank algorithm best ranks the features for best credibility assessment. The system was trained and tested on two fetched datasets [138].

A Multi-Stage Model [21] having: Relative Importance, Classification, and Opinion Mining Components. The system's Dataset was constructed using 1.2M Tweets of Topic: Iraq and Levant (ISIS) DAISH. Only 1000 tweets of 700 Users were tagged to train the Naïve Bayes classifier with Relative Feature Importance implemented over user and tweet level features. First of all complete User's Sentimental and Credible Tweets Ratio was computed, then Tweet's Credibility probability value predicted using a trained classifier and finally, both values are combined as weighted credibility score.

Total 2000 trendy tweets of 10 topics posted in japan were annotated through four questions and trained a Random Forest classifier. Four distinct features: tweet topic, user topic, user's expertness, and bias are additionally assessed. Tweet topic and user topic features were extracted from LDA and concluded that topical features improve credibility assessment [139].

Following are some serious observations, first: it has been observed across all automatic credibility assessment systems of any type (e.g.: Source based, Feature based, Graph based, and Hybrid) and even in explanatory studies, that majority of these studies get their dataset labeled either considering that: post seems 'informative/newsworthy' or 'trustworthy/truthful' to the evaluators. Only couple of studies considered Real and Fake news from authentic sources and get their dataset labeled on authentic basis rather than on evaluator's perception. Second: Many important aspects regarding evaluation criteria discussed in [217] are also fully ignored. Third: another important observation regarding every research study that they just consider news event for credibility, any other piece of information is not even considered for credibility assessment, though information credibility exist in every piece of information.

#### E. GRAPH BASED MODELS FOR CREDIBILITY:

Such studies of Source Credibility type, are classified in this category which uses Social Network Analysis (SNA)/ Graph-based models [218] by utilizing friendship (follower/following) network, user's tweet/retweet propagation network, etc. The majority of Source Credibility studies are graph-based (see section VI-C). An academic research (TURank) [128] is discussed as an example case which is classified as source credibility using the graph-based method. In this study, the original Twitter information network flow is

**TABLE 8.** Feature Based & Hybrid Microblog Credibility Assessment Models Classification ( table 2 of 2). Each study is comprehensively presented across following important attributes. They are salient qualitative dimensions of these research studies which should be known for efficient exploration of the research area.

Paper	Algo	Learning Type	Approach	Features Level	Dataset	Outcome	Label Method	Focus	Product	Distinct Attributes	Category
[137]	Random Forest	Classification	Feature Based	User, Tweets	7000 Tweets on Nature Environment Preservation, 100 Terms used, Tagged 1206 Tweets	Credible	Manual	Trying to make time efficient Twitter Plugin	Twitter Plugin	Reconciliation System for Tagging Evaluation & Retagging	ML
[21]	Naïve Bayes + Relative Feature Importance	Classification	Feature Based & Weighted Score	User+ History, Tweets	1.2M Tweets of Topic Iraq & Levant (ISIS) DAISH, Tagged 1000 of 700 Users	Credible	Experts	1.Complete User Level Sentimental & Credible Tweets Ratio is computed. 2. Tweet's Credibility Probability value predicted. 3. Weighted 1 & 2	No	Multi- Stage Model having: Relative Importance, Classification, Opinion Mining Components	ML
[31]	Random Forest (10 fold CV)	Supervised & Unsupervised	Feature Based & Graph Based	Topic+ Aggregated, User, Tweet, Directed Weighted Graph of: User, Tweet, Topic	Turkey's 25 Trendy Topics & 100 Tweets each & also Tagged	Labels: Important, Newsworthily, Correct	Crowd Worker	1.Different Models trained separate for user, topic, tweet and combined. 2.Labels are converted in scores. 3. Authority transfer 4. Threshold is used for getting labels	No	1. Authority Transfer Model by means of equations. 2. Weighted Directed Network of user-tweet-topic.	Hybrid
[140]	SVM	Classification & Unsupervised	Feature Based & Graph Based	User, Tweets, Events+ Aggregated, Weighted Directed Hierarchical Net. of: Event, Sub-event, Msg.	Cina Weibo's Two Datasets: 1.SW2013 (Topic Independent) 18 Fake, 171Real News (79K Tweets, 63K Users) 2.SW-MH370 (Topic Dependent) 32 Rumor, 103 Real (31K Tweets, 24K Users)	Credible	Fake, Rumor, & Real News are Collected from Authentic Sources	1.Focus on News Credibility. 2. Sub-event & Message Layers has Inter & Intra Layer Links. 3. Event & Sub-event Credibility Initial Scores are Calculated By Avg. of Message	No	Graph Optimization Method Used for Proposed Model Named NewsCP	Hybrid
[141]	Decision Tree(J48)-D2010 & KNN-D2011	Classification & Unsupervised	Feature Based & Graph Based	User, Tweet, Event, Directed Weighted Network of: Event, Tweets, Users	1. D2010-47K Users, 76K Tweets, 2K Topics, 207 Events. 2. D2011-9K Users, 76K Tweets, 2K Topics, 250 Events. (Events are tagged by 10 sample tweets)	Credible	Authors	1. Focus on Event Credibility. 2. Event & Tweet Implications computed as +ve/-ve weights within layers. 3. Event & Tweet layers have inter & Intra links	No	1. Initially used PageRank like Algo. as BasicCA. 2. Graph Optimization used to enhance results as: EventOptCA	Hybrid

used to find the authoritative user. The philosophy of TURank<sub>i</sub><sup>198</sup> says that: user becomes more authoritative when followed by<sub>f</sub><sup>199</sup> another authoritative user. Likewise, tweets become more im-<sub>i</sub><sup>200</sup>portant when retweeted and it also affects its user's authority.<sub>i</sub><sup>201</sup> Therefore types of such authority transfer in TURank are:<sub>i</sub><sup>202</sup> user-user, tweet-tweet, tweet-user, and user-tweet.<sub>i</sub><sup>203</sup>

Other graph-based models are intentionally not discussed for<sub>f</sub><sup>204</sup> these few reasons. 1. They are too many in quantity because<sub>e</sub><sup>205</sup> a majority of source credibility studies are all graph-based.<sub>i</sub><sup>206</sup> 2. There are surveys available for these graph-based models<sub>s</sub><sup>207</sup> which are explicitly discussed under source credibility. 3.<sub>i</sub><sup>208</sup> Important concepts and techniques are completely covered<sub>e</sub><sup>209</sup> in the other two types of models like feature-based models and hybrid models.<sub>i</sub><sup>210</sup>

#### F. HYBRID MODELS FOR CREDIBILITY:

Hybrid models combine the strength of both feature-based<sub>e</sub><sup>213</sup> and graph-based models, therefore a much better approach<sub>e</sub><sup>214</sup> has resulted in very few shortcomings. It is commonly ob-<sub>e</sub><sup>215</sup>served that studies in this area initially exploit feature-based<sub>e</sub><sup>216</sup> models to get User, Tweet, etc. seed scores which become<sub>e</sub><sup>217</sup> nodes of some user-defined network. Afterward, the network<sub>i</sub><sup>218</sup> of such entities like Topic, Tweets, Users, or Events, having<sub>e</sub><sup>219</sup> inter and intralayer-directed links with signed weights, are<sub>e</sub><sup>220</sup> made. Event/Topic initial scores may be generated through<sub>e</sub><sup>221</sup> aggregated values of their decedents. Finally, graph-based or<sub>f</sub><sup>222</sup> graph optimization methods are used for score convergence,<sub>i</sub><sup>223</sup> and some thresholds are used for credibility prediction. It<sub>f</sub><sup>224</sup> is explored that simply linking entities as a network enable<sub>e</sub><sup>225</sup> hybrid models to best exploit implicit entity relations.<sub>i</sub><sup>226</sup>

Following are the studies categorized as hybrid models for<sub>f</sub><sup>227</sup> credibility. In the study, [31] a total of 41 features for Topic,<sub>i</sub><sup>228</sup> Tweet, and User are used for learning and score generation.<sub>i</sub><sup>229</sup> As each tweet refers to the user as well as the topic, therefore<sub>e</sub><sup>230</sup> initial score is used in authority transfer for calculating the<sub>e</sub><sup>231</sup> credibility of each tweet. Dataset was generated through 25<sub>i</sub><sup>232</sup> trending topics of Turkey having 100 tweets in each.<sub>i</sub><sup>233</sup>

Another hybrid approach is used in [140]. Two Datasets<sub>s</sub><sup>234</sup> (topic dependent and independent) were used. Both extracted<sub>e</sub><sup>235</sup> from Cina Weibo's messages having Rumors, Fake, and Real News which were selected from authentic sources. The<sub>e</sub><sup>236</sup> SVM classifier was trained first on the user, tweet, and event<sub>e</sub><sup>237</sup> (aggregated) level features, and then a weighted directed<sub>e</sub><sup>238</sup> hierarchical network of entities as Event, Sub-event, and<sub>e</sub><sup>239</sup> Messages was constructed with inter and intralayer links.<sub>i</sub><sup>240</sup> Inter-layer links represent explicit relations between network<sub>i</sub><sup>241</sup> entities. Messages' initial credibility scores were generated<sub>e</sub><sup>242</sup> by a trained SVM classifier and then Event and Sub-event<sub>e</sub><sup>243</sup> credibility initial scores are calculated by respective aver-<sub>e</sub><sup>244</sup>ages. Finally, the graph optimization method was used for<sub>f</sub><sup>245</sup> the proposed model named NewsCP.<sub>i</sub><sup>246</sup>

In [141] two datasets having 76K Tweets and 2K topics<sub>s</sub><sup>247</sup> each, of 457 total Events, all were tagged with 10 sample<sub>e</sub><sup>248</sup> tweets. First of all two separate classifiers: Decision Tree<sub>e</sub><sup>249</sup> (J48) and KNN were trained for each dataset, on User,<sub>i</sub><sup>250</sup> Tweets, and Event level features then a weighted directed<sub>e</sub><sup>251</sup>

network of entities having Event, Messages, and Users was constructed. Entities were linked with their explicit relations. Event and Tweet Implications are computed as positive/ negative weights within each respective layer for their intralayer links. Initial tweet scores were obtained from the respective classifier and then a PageRank-like algorithm named BasicCA was executed over the network. The final optimized results were obtained from Event Graph optimization-based algorithm named: EventOptCA.

All above hybrid credibility assessment studies are summarized in Table: 8 (see last three entries of the table), across different attributes.

#### VII. STANDARD CREDIBILITY DATASET:

It is extremely important to discuss that one more challenging issue which is unsolved. It is the absence of predefined credibility benchmarks and its related gold standard dataset. The difficulty of collecting a large amount of such data has not yet received the attention it deserves [29].

Though there are many Deception related (e.g.: fake news, Rumor, Hoax, Spam, etc.) datasets (e.g.: LIAR [219], Fake-NewsNet [220], BuzzFeedNews [221], DeClare [222], Fake-NewsAMT [223], Hoaxy [224], Kaggle's- BSDetector [225], SemEval Task8 [226], Rumors [227], etc.) [228] are available. Web site's contents related credibility dataset [186], Event Credibility dataset [164], Bot and Malicious Profiles Detection dataset [101] and similarly few other credibility related components datasets are also available.

We have developed Credibility Taxonomy in table: 1 and figure: 3, summarizing all above sections (3-7) and the detailed classified tables: 7 and 8 to summarize and categorize automatic credibility assessment approaches across various dimensions, for all feature-based/ML/IR and Hybrid models. Graph-based models are intentionally not included for few reasons: one they are too many in quantity, second there are surveys available for only source credibility, and last; important concepts and techniques are completely covered in the other two as well.

#### VIII. LITERATURE BASED IMPORTANT FEATURES:

It is very important to know that what features are being used in microblogs credibility assessment studies, throughout the literature. Therefore, in this section most common and important features are extracted without any specific consideration of type, and methodology used. In this research study, there were almost 50 papers which were focusing specifically on microblogs. These were all discussed under section VI: Information Credibility of Social Media and Microblogs. There are two components in every information shared at microblogs: Post and Poster. At poster level: it is found that user's followers and followings, number of posts, age of account were found dominating in many papers. Location, picture in profile, description in profile were moderately used. It can also be observed that in the same user object, that time zone and gender are not much used (see figure 4).

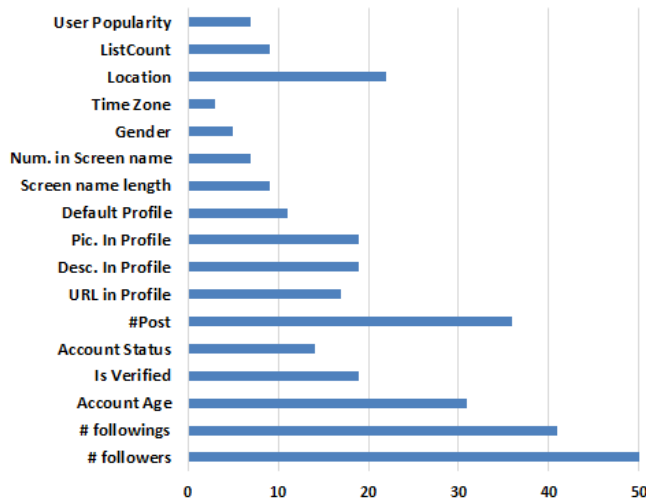


FIGURE 4. Mostly used User-related features in literature: in 50 papers.

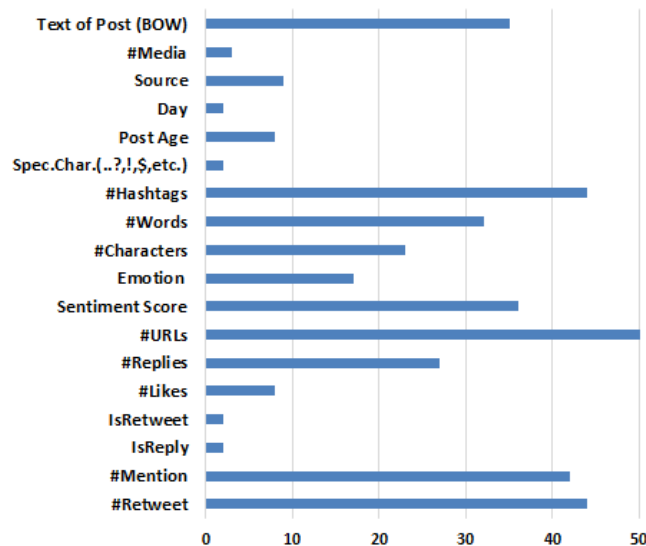


FIGURE 5. Mostly used Post-related features in literature: in 50 papers.

In post object: URL, retweet, hashtags, mentions are found strongly dominant in the majority of papers. Sentiment score and post content/text which was mostly used as bag of word (BOW) form, were also considered good features for assessing the microblogs' credibility. The number of words, number of characters, and number of replies are moderately used. Similarly, in post object few features like: number of media, isreply, isretweet, special characters, and day of the week are less utilized (see figure 5).

It is found that in addition to raw features (e.g.: #retweet, is\_reply, #mentions, #hash-tags, etc.) aggregated features and historic features performed better in assessing credibility [229].

The above most commonly used features are also adopted for the proposed framework's features presented in table 15.

## IX. FINDINGS AND DISCUSSIONS:

After a deep exploration of the literature and in-depth study of many automatic credibility assessment models (discussed from section VI-C to VI-F). These models or research studies are broadly categorized as following four types, to understand and briefly discuss their distinct features (see table 9), strengths, shortcomings (see table 10) and for recommendations (see table 11). Each category is briefly discussed under following respective headings and summarized in table: 9 as well:

**1. Feature Based - Tweet Credibility:** a large number of researches only extract features based on authors, contents of tweet, topic, and underlying network-related static features, e.g.: number of followers and followings, etc. (available in a tweet only) and apply Machine Learning models to identify credibility score or label. Such models completely ignore the influence of user's friends- network and post propagation networks, etc. On the other hand, they are also unaware that some very important credibility features like the number of retweets, likes, followers, etc. are generally inflated by malicious profiles/ bots, hence produce a completely false sense of credibility. Similar to all other categories they are also affected by the absence of post quality-related many credibility aspects (e.g.: Fake, Bias, Spam, Rumor, Smear Campaigns, Conspiracy, etc.) proposed in our credibility framework's table 14.

**2. Graph Based - User Credibility:** It is assumed in this category that if the post is authored or propagated through a highly authoritative or influential/central user then it is likely to be more credible. Many attempts are made using graph-based (un-supervised) methods to identify the user influence within the social network and then credibility is judged for the author's influence, and therefore infected with things like fake followers/ follower's fallacy, coordinated behavior, etc. These manipulations are mostly done by malicious profiles/bots. Focus is fully shifted towards the source of the message and therefore post itself is completely ignored and similarly, post-quality-related important credibility aspects discussed earlier are also ignored. Regarding automatic credibility assessment, we need to carefully set a threshold for identification of our credibility label, as data will be unlabeled for such un-supervised problems.

**3. Featured Based - Tweet + User Credibility:** It is the extension of 1st Category (named as Feature Based-Tweet Credibility), discussed earlier with additional focus on user-related credibility aspects. In addition to message/post credibility, by using different ways and means, the credibility of the user is also measured to label or score the final credibility. For example, assessing the credibility of the user, different historic or aggregated or weighted features could be used ( historic, aggregated features are discussed in sub-section VI-D). It is generally observed that; each message affects the credibility of its author and vice versa. It's an extremely

**TABLE 9.** Summary/ Important Aspects of Microblogs based Automatic Credibility Assessment Models categories (table 1 of 4). There were four types of Microblogs based Automatic Credibility Assessment Models (see sections: VI-C, VI-D, VI-E, VI-F).

1st Category of Research	2nd Category of Research	3rd Category of Research	4th Category of Research
<b>Feature Based – Tweet Credibility</b>	<b>Graph Based – User Credibility</b>	<b>Feature Based – Credibility: Tweet + User</b>	<b>Hybrid - Feature Based + Graph Based</b>
Using only or mostly Tweet level features, ML model is trained.	1. Using friends-following or User-tweet-RT Network, apply modified Page Rank like model or some form of Authority transfer, etc. for calculating highly influential/popular/reputed user as credible.	User level (Historic, Aggregated/ Weighted Features: Sentiments, Favorites, Mentions, Retweet, Listed, Friends NW influence, etc.) & Tweet level features are used to train ML model.	1. Initially feature based models are used to get user, tweet score, then network of entities (like: Topic, Tweets, Users, Events) having inter and intra layer links with weights, are made and finally some graph based methods used for score convergence.
	2. Ranking User.		2. Best exploits implicit entity relations.
	3. Un-supervised.		3. Much better approach with minor shortcomings.

**TABLE 10.** Shortcomings of Microblogs based Automatic Credibility Assessment Models categories (table 2 of 4).

1st Category of Research	2nd Category of Research	3rd Category of Research	4th Category of Research
<b>Feature Based – Tweet Credibility</b>	<b>Graph Based – User Credibility</b>	<b>Feature Based – Credibility: Tweet + User</b>	<b>Hybrid - Feature Based + Graph Based</b>
<b>Shortcomings</b>			
Credibility not captured.	Assumption that Post quality is just based on user.	Effectuated with followers fallacy, and bot manipulations like problems.	Assessing Event/Topic credibility is complex. Somehow tweet credibility.
Fakeness not assessed.	Post Quality is completely ignored.	Post Fakeness not assessed.	Respective Feature Based shortcomings are inherent.
Bot manipulations in RTs, Likes, Mentions, #Tags etc.	Bots/Cyborgs seems highly influenced here. Credible/Not-Credible not labeled	Bot manipulations in RTs, Likes, Mentions, #Tags etc.	Thresholds may be different for different nature of topics/events in real-time.

**TABLE 11.** Proposed or recommended features selected or considered under each research category of Microblogs based Automatic Credibility Assessment Models (table 3 of 4).

1st Category of Research	2nd Category of Research	3rd Category of Research	4th Category of Research
<b>Feature Based – Tweet Credibility</b>	<b>Graph Based – User Credibility</b>	<b>Feature Based – Credibility: Tweet + User</b>	<b>Hybrid - Feature Based + Graph Based</b>
<b>Proposed/Recommended Credibility System</b>			
<b>Hybrid (Feature Based + Graph Based) – Tweet Credibility Score: User + Tweet</b>			
Comprehensive Tweet level Features (in addition to others) are used to assess post quality rank through Learn to Rank Model.	Un-usual to this category, Graph Based models are applied at both User + Tweet levels. Graph Based models are applied at each tweet’s retweet network and user followers-following network.	Many User + Tweet level features, including Historic, Aggregated and simple features are considered.	Both User (Friends Network-Influence) and Tweet level (Retweet Network-Spread & Propagation) features having scores, which are used as features.  Finally all features including Network Scores and remaining normal User Features + Tweet Features are used to rank tweet, using Learn to Rank models.
	User influence score (using only trustworthy or Non-Bot followers-following network)		
	Tweet Spread, and Propagation scores (using retweet network) are also calculated.		

**TABLE 12.** Summary of completely distinct recommendations which are not considered in any of the four research categories presented in tables 9,10, 11 (table 4 of 4).

<b>Completely Distinct Considerations Recommended (not included in any research category)</b>
Identify if A/C behavior is like malicious Bot/Troll/Cyborg/Sybil, etc. (such A/C will be omitted from friends network for correct influence calculations)
Identify post fakeness (Post Level) and also update User’s fake producer counter (User Level)
Score of post is computed (using all User & Post Level features + actual retweet network’s propagation and spread measures + user rank over friend’s network)
Users features includes: Domain (area of expertise), correct influence calculated only over trustworthy friends network, etc.
User includes: Fake Produced % age, Spread Score & Propagation Score ( Avg. of Tweet Spread & Propagation Scores)

1323 important phenomenon but very rarely identified. We ob-1325 served that only one study tries to identify the credibility1326 score through identification of the topic of the tweet and then identify the number of topic-specific influential users

involved in re-tweeting and then determine the credibility score. One study identifies credible tweets only when if it remains original and then scores its source and then the tweet score is calculated but it's all about theoretical. One study proposed that if more authoritative/centric people are involved in retweeting then score of credibility is increased.

**4. Hybrid - Feature Based + Graph Based:** The modern method, which isn't sufficiently explored in studies till now, exploits the power of both Feature-Based and Graph-Based models, known as a hybrid. They attempt for Feature-Based models for initial credibility prediction of respective entities, for example, predict credibility of the tweet, user, topic etc. and then further boost the results through incorporating their scores/predictions to an interconnected network of participating entities like Post, Poster, Topic, and Event. There is an obvious observation, as we discussed in the above 3rd category, that each entity affects the credibility of others and gets affected, which means all are interdependent which is implicitly exploited through network models in hybrid settings. Despite the strength we discussed, they inherently suffer from some shortcomings of feature-based models as well. Few other shortcomings include difficulty in assessing real event credibility or topic credibility values which somehow primarily, tweet credibility again, e.g.: the credibility of a topic is computed through all their tweet credibility values. Once they are calculated, then they are again used in their interconnected network of participating entities, where these values are mostly amplified with some scalar effect. Another limitation is threshold settings which differ for the different domains (e.g.: politics, education, entertainment, sports, etc.).

**Shortcomings** (other than above categories): besides all above category-specific shortcomings there are some other extremely important shortcomings that are not discussed in any category because they don't fall in any category and they are also considered as our recommendations (see table 12). They are also discussed in section XII with other associated details and enlisted as following, like:

1. A very vital aspect that is completely ignored that the credibility of a message can't be determined without going into the underlying credible and trustworthy friend network, to measure the correct influence of the user. If malicious profiles exist in a friend's network then they must be omitted before examining the user rank/influence. Malicious profiles/bots identification and their rectification must be done for credibility assessment initialization, to prevent their serious manipulations at various places.
2. Chain of narrators is extremely important in assessing the message's credibility. Once a post is identified as fake then its producer must be penalized by incrementing its fake producer counter. Similarly, each fake propagator involved in post propagation within the post's chain of narrators must also be updated.

3. Credibility of the post must be calculated using a comprehensive list of features provided in table 15. This proposed list of features covers the majority of aspects like, post quality (which is ignored in the majority of studies, see figure 6 for a post-quality-related group of aspects), veracity and different forms of deception, hate speech, post spread and propagation, user's veracity, expertise, rank, and malicious profile identification. All of these features are extremely important for automatic credibility assessment, e.g.: the spread and propagation pattern of a message is an important feature for credibility assessment. Computing user's influence or rank on a followers-following network comprising of non-malicious users/profiles only. After computation of all such features, an appropriate Machine learning model could be trained over these features for score/rank prediction.

4. Two extremely important features which are fully ignored in credibility studies are user domain/topic-specific expertise and true user influence score computation without bot manipulations.

5. As it has been discussed earlier that many post-level features could compute user-level features. Therefore many user-level scores could easily be computed like, User's Avg. Post Credibility Score, User's Fake Post Produced %age, User's Fake Post Propagated %age, User's Spread and Propagation Score Avg., etc. Computation of all such scores at the user level will implicitly reduce the dissemination of low-credibility contents, over microblogs.

Detail recommendations are presented in section XII.

We have also presented a summary of the above observations in table 9 with their shortcomings in table 10 and our those recommendations which are based on already defined research categories, presented in the table: 11, whereas recommendations which are fully distinct or completely missing in all the categories are proposed in table 12.

## X. THEORY DRIVEN CREDIBILITY FRAMEWORK:

The framework has theoretical foundation. How the framework is driven and what are the basis of our proposed framework is presented as following:

1. Basic components (Levels, Dimension, Constructs) of credibility are identified through detailed literature exploration from different disciplines of credibility like physiology, communication, information sciences, etc. (see section III-A under heading 'Credibility Components', and table 2 and 3).
2. All credibility supported research studies were identified first, after detailed literature exploration, then each concerned research study is categorized and discussed under its respective construct. Example: Fake News Detection studies are categorized and discussed under Deception, Truthful constructs (see complete section V).
3. Necessary credibility components identified in step 1 and 2, are presented in the form of a framework, presenting their inter-relationships (see table 14).

**TABLE 13.** Economics, Social Sciences Basic Theories, and Credibility Studies Driven Credibility Framework Components.

Basic	Framework	Components	Theory	Description	Research Based Ref.
Post Related Theories	Contents	Quality	Information Manipulation Theory [230]	Too many or too few refers to deception	It is primary component so too many ref. are found, just few are: [136], [157], [38], [167], [231], [158]
			Reality Monitoring [232]	Real events are identified by sensory perceptual information	
			Four Factor Theory [233]	Emotion, arousal, thinking, and behavioral control are expressed differently in lies and truth.	
			Undeutsch Hypothesis [234]	Factual contents differ in quality and style from fallacy	
Source (User) Related Theories	Expertise	Community/ Peer Influence	Rare Behavior [18]	Unusual behavior than majority	Expertise: [235], [236]
			Synchronized Behavior [18]	All such user show/ follow the similar behavior patterns.	
			Coordinated Behavior [18]	All such chain of users are developed to perform some pre-defined task of their master.	
			Collective Behavior [18]	Actions performed by presence oriented mass (crowds, mobs, riots, cults)/ distance oriented (rumors, mass hysteria, moral panics, fads, crazes)	
			Social Identity Theory [237]	portion of an individual's self-concept derived from perceived membership in a relevant social group	Combined Expertise & Trustworthiness [43]–[46], [238], [38], [154], [239]–[241], [22], [231], [242], [243]
			Emperor's Dilemma [244]	Alternative possibility, that members of a group may enforce to act in ways that few if any group members actually want or need.	
			Normative Influence Theory [245]	People change to form a good impression and fear of embarrassment or to be liked or accepted by others	
			Availability Cascade [246]	Self-reinforcing process in which a collective belief gains more and more plausibility through its increasing repetition in public discourse within their social circles	
	Trustworthiness	Individual Influence	Overconfidence Effect [247]	One's subjective confidence in his or her judgments is reliably greater than the objective ones.	Trustworthiness: [36], [254]–[256]
			Illusion of Asymmetric Insight [248]	We understand others better than they understand themselves	
			Naïve Realism [249]	A believe that we see the world objectively, and people who disagree, must be irrational, or biased.	
			Selective Exposure [250]	Prefer information based on pre-existing attitude.	
			Confirmation Bias [251]	Trust information based on pre-existing beliefs.	
			Desirability Bias [252]	Accept information that please them.	
	Driven By Benefits	Community/ Peer Influence	Bandwagon Effect [253]	Do something because others are doing.	Trustworthiness: [36], [254]–[256]
			Conservative Bias 158	Revise one's belief insufficiently when presented with new evidence.	
			Validity Effect [257]	Believe that information is correct after repeated exposures.	
			Semmelweis Reflex [258]	When something contradicts with well established norms then reject such new evidences	
Attentional Bias [259]			failure to consider alternative possibilities when occupied with an existing train of thought		
Echo Chamber Effect [260]			Within a close system, belief are amplified by communication and repetition		
Contrast Effect [261]			When compression enhances differences.		
Prospect Theory [262]	People decide between alternatives like gains or losses, and just think in terms of expected utility rather than absolute outcomes.				
Optimism Bias [263]	Overestimate the probability of positives and underestimate the probability of negatives				



1438 **4.** To strengthen our framework components we identify<sup>1494</sup>  
1439 the basic theories of Economics and Social Sciences which<sup>1495</sup>  
1440 are supporting or leading towards individual framework<sup>1496</sup>  
1441 components (see table 13).<sup>1497</sup>

1442 **5.** To strengthen our framework components we identify<sup>1498</sup>  
1443 the basic credibility studies which are supporting or leading<sup>1499</sup>  
1444 towards individual framework components (see table 13).<sup>1500</sup>

1445  
1446 It has already been discussed that outcome of our study was<sup>1502</sup>  
1447 two-fold. Understanding the broader domain of credibility<sup>1503</sup>  
1448 with basic components identification (i.e.: levels, dimensions,<sup>1504</sup>  
1449 and constructs) and then the development of compatible<sup>1505</sup>  
1450 social media generic framework will be carried out. This<sup>1506</sup>  
1451 will further be transformed to microblog-specific imple-<sup>1507</sup>  
1452 mentation. Considering the first objective, a theory-driven<sup>1508</sup>  
1453 generic framework of social media is going to be identified<sup>1509</sup>  
1454 in this section, consisting of levels and dimensions. These<sup>1510</sup>  
1455 two components are completely generic to social media only.<sup>1511</sup>  
1456 Constructs must be carefully identified for both social media<sup>1512</sup>  
1457 and microblogs. Therefore they are identified in the next sec-<sup>1513</sup>  
1458 tion XI in addition to our microblog specific implementation<sup>1514</sup>  
1459 as our second objective fulfillment. The generic (levels and<sup>1515</sup>  
1460 dimensions) and specific (constructs) framework components<sup>1516</sup>  
1461 have already been identified in previous section III, under<sup>1517</sup>  
1462 the heading of 'Credibility Components', through strong and<sup>1518</sup>  
1463 detailed literature exploration.<sup>1519</sup>

1464  
1465 In addition to the literature explored in previous sections,<sup>1521</sup>  
1466 to form the strong basis of credibility framework. A compre-<sup>1522</sup>  
1467 hensive and dual study is also conducted as follows. Table<sup>1523</sup>  
1468 13 completely map our framework components (see first<sup>1524</sup>  
1469 merged column for framework components) with the fol-<sup>1525</sup>  
1470 lowing Social Sciences & Economics Theories (see second  
1471 column for these theories with short description) and then<sup>1526</sup>  
1472 with Credibility Studies in the last column (see research-<sup>1527</sup>  
1473 based references of these studies):<sup>1528</sup>

1474  
1475 **10.1. Social Sciences & Economics Theories Driven:** We<sup>1530</sup>  
1476 have surveyed many related basic behavioral and human<sup>1531</sup>  
1477 cognition theories defined across varied disciplines: like<sup>1532</sup>  
1478 economics and social science. Each theory with its short<sup>1533</sup>  
1479 description is presented in table 13. They provide important<sup>1534</sup>  
1480 guidelines for the required level ( post or poster) of credibility<sup>1535</sup>  
1481 and deception. Such theories simply lead towards building<sup>1536</sup>  
1482 efficient models of credibility identification or assessment.<sup>1537</sup>  
1483 High-level analysis of these selected theories resulted that<sup>1538</sup>  
1484 they are either related to the post itself or posters. Hence two<sup>1539</sup>  
1485 pillars or levels of credibility could be identified first which<sup>1540</sup>  
1486 are 'post' and 'poster'. Further considering the important<sup>1541</sup>  
1487 dimensions of credibility. These theories are also classified<sup>1542</sup>  
1488 under 'content quality' of post and two types of influence<sup>1543</sup>  
1489 (e.g.:community and individual) which directly affect either<sup>1544</sup>  
1490 'poster's expertise' or 'trustworthiness'. Some specific the-<sup>1545</sup>  
1491 ories are driven by benefits that affect the poster's trustwor-<sup>1546</sup>  
1492 thiness as well. Therefore three major dimensions are also<sup>1547</sup>  
1493 identified: content quality, expertise, and trustworthiness (see<sup>1548</sup>

table 13).

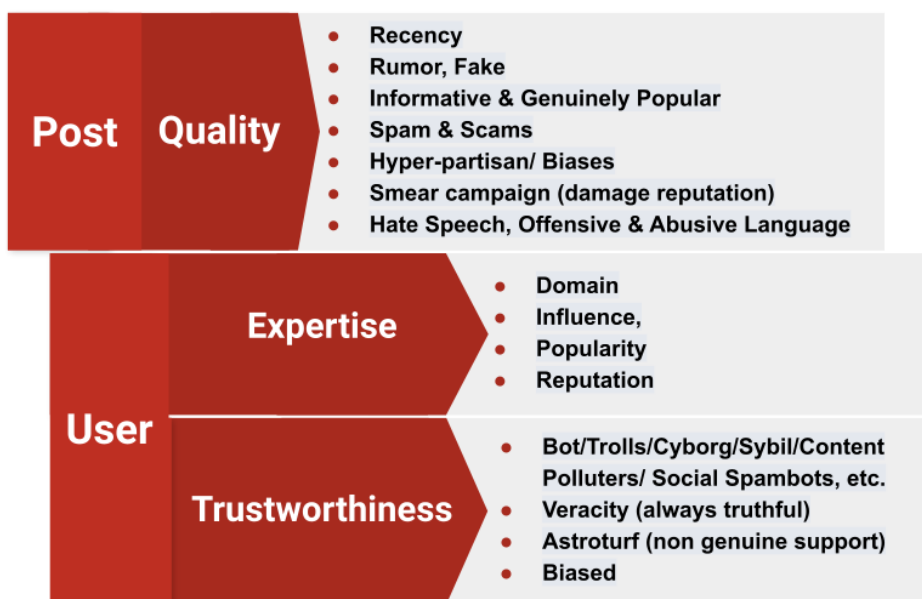
**10.2. Credibility Studies Driven:** In strong support of our framework components, we had already explored detailed credibility studies and credibility macro components (e.g. Levels, Dimensions and Constructs) had also been identified in section III (see table 2 and 3). It has also been discussed that considering social media credibility, only two main levels of credibility named message credibility and source credibility are feasible [42]. Regarding media credibility, it is also discussed earlier that in modern scenario medium is also replaced with source only [59]. In this case, the source has to be thoroughly examined including all chain of narrators involved in message propagation. Many leading credibility related research studies highlighted 'Trustworthiness' and 'Expertise' as major dimensions of source credibility (see last column named 'Research based Ref.' of table 13) and also table 2, 3. It could also be seen in table 13 that content quality is the most important and primary dimension of message/post (see table 2 references [142], [145], [157], etc. and table 13 references [38], [136], [158], [167], [231], etc.)

Identified credibility framework which is completely generic to social media, is theory-driven. The framework is fully supported through social sciences & economics theories and credibility-related research studies. Complete mapping of these theories and important research studies are all provided in a single comprehensive table 13. Considering the primary objective of the study, the extract of credibility framework is theory and research studies driven. High level credibility framework picture is further presented in figure:6, which will be completely understood after section XI.

## XI. PROPOSED SOCIAL MEDIA CREDIBILITY FRAMEWORK:

The framework was identified through supported theories in the previous section. That theory-driven framework is presented with further necessary details, in this section:

**11.1. Motivation and Objective:** Determining the credibility of information in microblogs is becoming one of the most challenging issues day by day and still, unresolved [20]–[22]. Even though it has been studied much, since last many years. It is observed that too much work is done on theoretical or conceptual aspects of credibility in other related fields but they are not properly considered in microblogs related automatic credibility assessment studies conducted in computer science. These theoretical or conceptual aspects of credibility are mostly studied in psychology, communication and information sciences, where as microblogs related automatic credibility assessment studies are done in computer sciences fields. Unfortunately, no work had been done on mapping these general constructs of credibility for microblogs, which should be considered minimally when developing the respective system to assess the credibility. Due to being multi-perspective nature, the diversity in the definition and percep-



**FIGURE 6.** Generic Social Media Credibility Framework's High-level Component Diagram: Constructs are intentionally omitted from the picture for simplicity and understanding.

**TABLE 14.** Proposed High-level Generic Credibility Framework for Social Media: presenting relationships between credibility Levels, Dimensions, Constructs, and Aspects.

	Dimensions (Level)	Constructs: Aspects	Descriptions
<b>High-level Credibility Framework</b>	<b>Quality (Post)</b>	<b>Deception, Truthful:</b> Rumor and Fake	Misinformation, Disinformation, Hoaxes, etc.
		<b>Uniqueness/Completeness:</b> Informative. <b>Popularity:</b> Genuinely Popular <b>Recency:</b> Recency	General quality related attributes (Informative, Recent, etc.). Popularity must be clean from Bot manipulations.
		<b>Deception, Truthful:</b> Spam & Scams	Phishing, click-bait, Political Astroturf Meme, Fake Reviews, etc.
		<b>Unbiased:</b> Hyper-partisan/ Biasness	Polarization, etc.
		<b>Deception, Truthful:</b> Smear campaign (damage reputation)	Satire, Meme, Propaganda, Conspiracy, etc.
		<b>Plausibility:</b> Hate Speech, Offensive & Abusive Language	All types of Hates: Ethnic based, Xenophobia, Islamophobia, racism, misogyny, etc.
	<b>Expertise (User)</b>	<b>Competence/Topic/Specificity:</b> Domain	Top three areas in which he message, e.g: Politics, Sports, Health, etc.
		<b>Authority:</b> Influence. <b>Popularity:</b> Popularity. <b>Competence:</b> Reputation	Different measures of Centrality, Authority Transfer, User Defined Ratios & Influence
	<b>Trustworthiness (User)</b>	<b>Deception, Truthful:</b> Bot/Trolls/Cyborg/Sybil/Content Polluters/Social Spambots, etc.	Fake A/Cs, Non Human Behavior, etc.
		<b>Truthful:</b> Veracity (always truthful)	Not fake and rumor producer/propagator, don't like them as well.
		<b>Deception:</b> Astroturf (non genuine support)	Followers fallacy, Bot Nets, Troll Factories/ Troll Farm, Link Farming, etc.
		<b>Unbiased:</b> Biased	If greater no of Hyper-partisan posts found

tion of credibility reflects different viewpoints in different work studies. Some studies consider 'Relevance' as the criterion of being credible. Some assume 'Reputation' as the major driver of credibility, whereas the majority only stick that 'Fake' identification is credibility identification. It is also perceived by researchers, that 'Ranking' concerning author influence and topic expertise are strongly treated as credibility ranking. The majority of studies exploit 'Informativeness' as a credibility indicator. Few found examining 'Trust' level as true credibility judgment. It is observed and quite evident in many research studies as well, that the credibility notion needs to be standardized because each one of them only covers some aspect of credibility, and the majority are left undiscovered.

One important objective of the study was to fill the specified gap and propose a theoretical framework with a similar approach followed in many similar studies like [51], [59], [144], [157], [166].

**11.2. Findings:** Investigating and exploring the credibility studies found in different fields, like psychology, communication, and information sciences, etc. identified extremely important credibility constructs under the dimensions and levels. There were some critical constructs also identified which were completely missing in many credibility assessment studies. Therefore challenge of credibility assessment was unresolved. Many research studies are now considered under these constructs. Following are some examples studies considered under their respective constructs: Hyperpartisan, Hate speech & offensive language, and Smear campaign which are considered under post quality's constructs Bias/Objectivity, Plausibility, Deception/Truthful respectively. Similarly some other studies like malicious profiles (bots, cyborgs, Sybils, etc), and astroturf (non-genuine support) which are considered under user trustworthiness constructs Deception/Truthful, Truthful respectively (see table 14).

It is also discovered that credibility is composed of many constructs, which are identified in section III-C2. All these constructs must be considered in assessment instead of considering only one or two. Majority of earlier credibility assessment studies only consider one or two constructs, like relevance, deception, truthful, popularity. Only these some construct were mostly considered in majority of the studies in isolated manner and remaining all were ignored.

**11.3. Constructs for Social Media and Microblogs:** The proposed framework is comprised of specific credibility levels, dimensions, and constructs and simply presenting their relationships. Credibility levels and dimensions identified in section III are general to social media credibility and therefore could serve as standard social media credibility framework, regardless of a very specific information object, whereas constructs will be information object-specific means both social media and microblogs specific [68]. Therefore in this section, such important constructs will be identified

and then a proposed social media framework will be presented. Distinct social media and microblogs characteristics are discussed in section VI. Important list of constructs are selected from table 2 and 3 considering the specified characteristics and presented as following. The same list of constructs was already shortlisted in section III-C2 with the same preferences.

The list is presented again for easy reference. These constructs together with associated aspects are also shown in the table 14, presenting high-level credibility framework:

1. Recency, 2. Truthful, 3. Deception, 4. Topic, 5. Specificity, 6. Unbiased/Objectivity, 7. Popularity, 8. Plausibility, 9. Authority/Influence, 10. Competence/ Reputation, 11. Uniqueness/ Completeness, etc.

Finally to complete our proposed generic social media-based credibility framework. Specific aspects/ characteristics elaborating each construct are also presented.

Considering third/ second last column of table 14. All constructs are specified in bold and aspects are written adjacent to the constructs. For example considering the first line, Deception, Truthful (constructs): Rumor and Fake (aspects). It simply means that 'Deception, Truthful' constructs could be implemented through 'Rumor' detection and 'Fake' detection. These 'Rumor' and 'Fake' are aspects, which need to be implemented for fulfilling respective constructs (e.g.: Deception, Truthful). All other remaining aspects are specified under their constructs, in the same manner.

The framework presents all components like Levels, Dimensions, Constructs, and related Aspects (see the complete framework in the table: 14). This framework will be further transformed for microblogs using microblog specific features, in the last segment of this section.

**11.3. Overview of Proposed Framework:** After conducting detailed and organized literature exploration it is proposed, that true Credibility is measured through narrator (user level) and their narration (post level) both (see figure 6). Narrator assessment may be done on its 'Expertise' and 'Trustworthiness' (see dimensions of the user), which are further assessed on multiple bases (see aspects, e.g.: domain, influence, popularity, reputation under 'expertise' dimension of the user). The narrator's 'expertise' could be judged through its genuine 'influence' based only on trustworthy social network context, level of expertise with relevant 'topic/domain', together with his/her 'popularity' and good 'reputation'. The narrator's trustworthiness could be assessed through the following aspects: the narrator should always be 'truthful', must not be 'biased'. The narrator should not behave like malicious profiles (e.g.: Bot/ troll/ cyborg/content polluter, etc.), etc. Similarly, narration may also be assessed on its 'Quality' (see the dimension of post). Quality may have different bases for assessment (see aspects, e.g.: recency, fake & rumor, hate speech, offensive & abusive language, biasness, informative, popular, etc. under the quality dimension of post). The quality of the post could be judged through different aspects: post

1661 'truthfulness', level of 'informativeness', and 'popularity'<sup>1717</sup>  
 1662 Post must also be clear from hate speech, and biasness, etc.<sup>1718</sup>  
 1663 (see figure 6).<sup>1719</sup>

1664 An effort is being made to present a proposed generic credi-<sup>1720</sup>  
 1665 bility framework (see table: 14) for social media. Comprising<sup>1721</sup>  
 1666 the levels (see column II - e.g.: Post, User) at which the<sup>1722</sup>  
 1667 credibility should be assessed together with respective di-<sup>1723</sup>  
 1668 mensions which completely adhere to the credibility related<sup>1724</sup>  
 1669 research studies and theories (see same column II – 1. Post;<sup>1725</sup>  
 1670 Quality. 2. User: Expertise, Trustworthiness) which need to<sup>1726</sup>  
 1671 be addressed. Finally what aspects/attributes under specified<sup>1727</sup>  
 1672 constructs (see column III – 1. Post Quality: Fake, Spam,<sup>1728</sup>  
 1673 Hyper-partisan, etc. 2. User Expertise: Domain, Influence,<sup>1729</sup>  
 1674 3. User Trustworthiness: Bot, Veracity, Biased, etc.) are<sup>1730</sup>  
 1675 comprising each construct under the dimensions. Last col-<sup>1731</sup>  
 1676 umn (see column IV) of our framework presents related<sup>1732</sup>  
 1677 or similar attributes which will be automatically covered if<sup>1733</sup>  
 1678 someone just considers the main aspects/attributes presented<sup>1734</sup>  
 1679 in column III. It could be noticed that the comprehensive set  
 1680 of aspects/attributes have mostly resulted from a thorough<sup>1735</sup>  
 1681 study of a large set of supported researches presented in<sup>1736</sup>  
 1682 section V. High-level credibility framework picture is also<sup>1737</sup>  
 1683 presented in figure:6. Important terms are also defined in the<sup>1738</sup>  
 1684 appendix section of the paper for clarity and understanding.<sup>1739</sup>

1685 <sup>1740</sup>

### 1686 **11.5. Framework Mapping to Microblog's Features**<sup>1741</sup>

1687 Considering the second objective of the study: the social<sup>1742</sup>  
 1688 media generic framework will be transformed to microblog-<sup>1743</sup>  
 1689 specific implementation through microblog's specific fea-<sup>1744</sup>  
 1690 tures. Therefore after presenting the most important base-<sup>1745</sup>  
 1691 line of our work as Proposed Social Media Credibility<sup>1746</sup>  
 1692 Framework. We are now presenting in the table: 15, that<sup>1747</sup>  
 1693 how each aspect/attribute of our social media credibility<sup>1748</sup>  
 1694 framework could be implemented over microblogs, through<sup>1749</sup>  
 1695 our proposed list of sample features. These features have<sup>1750</sup>  
 1696 mostly resulted from a detailed study of researches presented<sup>1751</sup>  
 1697 in section V and, section VI. Each feature is then justified<sup>1752</sup>  
 1698 by appropriate reference of research (covering a wide range<sup>1753</sup>  
 1699 of literature review, two complete sections of the study,<sup>1754</sup>  
 1700 section V and, section VI), together with its significance and<sup>1755</sup>  
 1701 judgment.<sup>1756</sup>

1702 The proposed list of sample features are furnished with<sup>1757</sup>  
 1703 two different levels (e.g.: User-Level, Post-Level), network<sup>1758</sup>  
 1704 features (e.g.: Friends network's Influence or Rank, Retweet<sup>1759</sup>  
 1705 network's Spread and Propagation), aggregated features<sup>1760</sup>  
 1706 (e.g.: Reciprocity, Reputation, etc.), and historic features at<sup>1761</sup>  
 1707 user level (e.g.: Domain, Veracity, Biased, etc.).<sup>1762</sup>

1708 <sup>1763</sup>

1709 Features presented in table 15 have varying levels of com-<sup>1764</sup>  
 1710 plexity. Few features are very simple and they are known<sup>1765</sup>  
 1711 as raw features, e.g.: number of followers, number of fol-<sup>1766</sup>  
 1712 lowings, age of account, is-verified, number of posts, URL<sup>1767</sup>  
 1713 in profile, description in profile, etc. Few features will<sup>1768</sup>  
 1714 be computed either through a separately trained machine<sup>1769</sup>  
 1715 learning system or by putting some extra effort, like the use<sup>1770</sup>  
 1716 of some lexicon, dictionary, etc. Examples of such features<sup>1771</sup>

are, Bot/Cyborg Likelihood Score, Hate-speech (Y/N), Abu-  
 sive Language (Y/N), Sentiment Score, Emotion Valance-  
 Arousal- Dominance (VAD) Score, Bias (Y/N), Fake (Y/N),  
 Topic of the Post (e.g.: Politics, Sports, Education, Social  
 Issues, etc.), Psycho-linguistic features calculated through  
 Linguistic Inquiry and Word Count (LIWC) lexicon with fol-  
 lowing categories: Informality, Cognitive Process, Perceptual  
 Process, and Diversity. Some features could be computed by  
 calling API, e.g.: Web Of Trust (WOT) Score, Informative:  
 Alexa Rank, Likes or Dislikes of YouTube Videos, and  
 Ground Truth Labels for the URL's found in the post. Some  
 features could be computed through standard libraries or self-  
 made programs. This list of features is: 'User Ranks' which  
 could be computed using page rank or modified page rank-  
 like algorithms, or different centrality measures of social  
 network analysis, etc. Other such features are Spread and  
 Propagation features of the post's re-tweet network which  
 could be computed using tree libraries.

## 1680 **XII. RECOMMENDATIONS:**

For better understand-ability, this section will present all the  
 recommendations as a blueprint, sketch, or glimpse of the  
 real system as if the system should look like this. Our basic  
 recommendations are presented under two tables. Table 11  
 presenting category-specific common properties which are  
 also found or picked in our proposed solution, whereas table  
 12 presets completely distinct and new properties which are  
 unique to our recommended solution.

The proposed solution is overcoming all identified short-  
 comings and further strengthening itself with extra proposed  
 features.

**12.1. Guided Data Tagging:** Data tagging is most important  
 for automatic credibility assessment systems. Following are  
 few serious issues found in these studies. Those are addressed  
 as following:

The required level of reliability needed in labeling the  
 credibility dataset would require a completely different pro-  
 cess. Data could not be tagged only based on the evalua-  
 tor's/expert's perception about the post. It is very challenging  
 to correctly label such multi-perspective data without dis-  
 covering hidden facts about the post. Data will be tagged  
 in a completely guided environment. Each post will be  
 tagged after various flags indicated by the variety of available  
 tools. All aspects of credibility must also be considered.  
 Expert/ evaluator will be indicated about poster's likelihood  
 score of the malicious profile, top 3 domains of the poster,  
 Avg. number of malicious profiles found in poster's friend  
 network, post's WOT score, Alexa rank, Ground Truth labels,  
 etc. if URL is found in the post, etc.

During tagging/scoring, all aspects of credibility must be  
 examined instead of only a few aspects which are mostly  
 examined in most of the studies. The majority of studies  
 either consider only fake/real as credibility. Some consider  
 that only popular, topic expert is representative of credibility,  
 etc.

**TABLE 15.** Implementing Social Media Credibility Framework to Microblog's: following are mapping of social media credibility framework's aspects to proposed microblog's sample features. Transforming the generic framework to microblogs specific implementation, just need the generic aspects to be transformed to microblog's features.

S.No.	Feature Name	Feature Level	Cred. Framework Aspects	Reference/Reason
1	User Ranks: Influence, other Centrality Scores, etc.	User Level (Friends Network)	Expertise, Quality	Measure of user influence and rank [34]
2	No. of followers	User Level	Bots, Expertise, Trustworthiness, Fake	Too few and too many: less expertise and trustworthiness. Less gap b/w 2 and 3: high competence, Ratio determines nature of A/C e.g.: broadcast, etc. Too many 2 and 3: Bot [121], High rate of friend/followers: Fake post producer [264], significant no of connections: active user [134]
3	No. of friends			
4	Age of a/c		Trustworthiness, Veracity, Expertise	Old A/C: produce less misinformation [264] and more trustworthy, Expertise, Competence [123] and New A/C: produce more misinformation and less trustworthy [264].
5	IsVarified and Protected		Bots, Fake	Verified a/c means real a/c notbot and Fake post producer [264].
6	No. of Tot.Posts		Trustworthiness, Expertise	High no of posts: credible post producer, active user [134]. User posting behavior:tweets/re-tweets [134]
7	URL in Profile (Y/N)		Trustworthiness	User perception based features visible at a glance, if yes then user perceive as credible [119]
8	Desc in Profile, Pic (Y/N)			
9	Bot/Cyborg Likelihood		Bot, Misinformation	Covers many aspects of credibility [23]. Bot: 0 or Very Less [265]. Bot:0, Human:1; Celebrities and popular org: high, more followers than followings. Bots: More followings than followers [98] Bot: Low, Human: High [98] Mostly non active , new user uses default [265].
10	List Count			
11	Reputation: Followers/Followers + Followings			
12	Reciprocity: fraction of friends who are also followers (overlap)			
13	Default Profile			
14	Domain (Top 3 domains extracted from post of user having topics)		Expertise, Quality	Once expert's domain and tweet topic is matched, fully reflects credibility [139].
15	Hate Speech, Abusive and offensive Language. (Y/N)		Tweet / Post Level	Hate, Quality, Smear
16	Get Ground Truth Labels for each URL in the Post.	Fake, Satire, Bias, Hate, Rumor, Spam Conspiracy		Varying level of Reliability and Bias labeling, URL's could be used for post identification as: Fake, Satire, Extreme Bias, Conspiracy, Rumor, Click-bait, Hate Group, Junk Science, etc. [183], [266]–[269]
17	Network: #Retweet	Quality, Fake, Rumor		Popularity, symbol of quality, msg endorsement [134], [138]. Considered important [30]. Fake has high retweets. [3], [84]
18	#mentions	Quality, Spam		Considered very important feature [30]. Too many mentions low credibility [123], in emergency also [132]
19	IsReply	Quality		One of some user perception based features visible at a glance, if yes: seems credible [119], it shows that User listen,agree/disagree and validate [264]
20	IsRetweet			
21	No. of Likes	Quality, Fake		Treated as good reputation [138]. Real news has more likes where as Fake has less [84].
22	No. of Replies	Fake, Bot		Bot: Very Less [270], Fake Post: High [20], Less [84]
23	Links: No. of URL	Fake		URL presence: High Credible [30], [134], Fake posts: large no of URLs [264]
24	WOT Score for URLs	Fake, Spam		Site reputation Score: Low score bad reputation [136] and spam, etc.Internet Trust Tool [271]

1772 Evaluators must be given clear guidelines for tagging, like 782  
 1773 what will be the credibility label/score/rank if the post 783  
 1774 is posted by a topic expert and the topic of the post is 784  
 1775 completely matched with the expertise of the poster. What 785  
 1776 label/score will be assigned to a post that is fake and posted 786  
 1777 by a malicious profile, etc. What about the post that has 787  
 1778 extreme bias and suffering from hate speech with abusive 788  
 1779 language. 1789  
 1780 In the presence of such indicators with clear guidelines post 790  
 1781 will be ranked/labeled finally by the expert/evaluator. 1791

**12.2. Hybrid System - Graph Based + Feature Based:** It is supposed to be a Hybrid system of a different kind. In our proposed solution: The graph-Based method will be executed first on two network-based features. There are two distinct sets of network features based on retweet network, and friends network, presented in Table: 15, at no 1, 46, and 47. Using friends network, where malicious profiles/bots will be eliminated, and the influence scores for each user will be calculated and saved as User level feature (see feature no: 1 in ta-

S.No.	Feature Name	Feature Level	Cred. Framework Aspects	Reference/Reason
25	Likes/Dislikes (if YouTube Video(s)), etc.	Tweet/ Post Level	Quality	High values good reputation and credibility
26	Psycho-linguistic (Informality): No. of Swear words/ Netspeak/ Assent/ Non-fluencies/ Fillers/ Typos		Fake, Quality, Spam	Psycho-linguistic LIWC [272](Informality) features. In news: Non Fake [134], [273], Identify type of tweet: Non News. Presence shows bad quality, Presence: Non Spam
27	No. of Self Words(i,my,mine)			Word like "I saw" more credible, Identify tweet type: Non News, Non Spam
28	Pronoun (1st, 2nd, 3rd Person) Present (y/n)			Identify tweet type: Non News, Non Spam
29	Sentiments: Sub/ Obj Score, etc.		Quality, Bias, Fake	Negative Sentiments are more credible in news. Generally either positive/ neutral is credible [134]. Real News: High ratio of neutral replies and Fake: High -ve replies. [84]. Bias Language Corpus [274]
30	Emotion VAD Scores			High Bias Language: High Fake.
31	Language Bias			
32	Text: Length			
33	No. of words		Quality, Fake	More length : more credible [30], [112], [134]
34	Fraction of upper case letters		Spam	High fraction leads to spam [164]
35	No. of Hashtags		Spam	3 or more are considered as spam [164]
36	?, !, Stock Symbol (\$) (Y/N). Contain multiple ?, ! (Y/N).		Fake, Quality, Spam	Identify tweet type: Non News, Non Spam (completely varying behavior in different aspects)
37	Smile icon, frown icon (:), etc. (y/n)			
38	MetaData: Age(sec)			
39	Day of the Week		Quality	Capture all time dependent aspects
40	Source (API 3rd Party, Un-Reg., mob, web)	Bot, Trustworthiness	Human: Web/Mob. Bot: API 3rd Party [98]. Source as Mobile is more credible.	
41	IsGeo-Coordinates (Y/N), etc.	Fake, Quality	Represent Location: More credible [136]	
42	Fake: Yes/No	Fake, Misinformation	Used for assessing truthiness, controls misinformation. 200 Fact Checking Web Sites [275]	
43	Topic: Politics, Health, Sports, Education, etc.	Expertise, Quality, Misinformation	If user's domain and tweet topic is matched, fully reflects credibility [139]. Some Topics are less credible [31], [131]. Misinformation is more diffused in some topics [179], [276].	
44	Informative: Alexa Rank	Quality	High Rank means informative and credible.	
45	Psycho-linguistic/LIWC: Cognitive Process, Perceptual Process and Diversity	Fake, Misinformation, Quality	Different classes of attributes are identified in LIWC [272] to identify Fake [273] .	
46	Spread: Level No., No. of RTs at each level (apply spread model)	Tweet/ Post Level (Retweet Network)	Fake, Misinformation	High spread and propagation lies in fake news, misinformation, etc. [20], [277], [278]. Very specific patterns are found in majority of Misinformation type contents within the Retweet Network [102], [175].
47	Propagation: Root Degree, Max Subtree, Avg. Subtree, Tree Max Degree and Avg. Degree (excluding root), Tree Max Depth, Avg. Depth			

ble 15). It is an extremely important aspect that is completely ignored that the credibility of a message can't be determined without going into the underlying credible and trustworthy friend's network, to measure the correct influence of the user. If malicious profiles exist in the friend's network then they must be omitted before examining the user rank/influence. Malicious profiles/bots identification and their rectification must be done before the credibility assessment initialization to prevent their serious manipulations at various places. Likewise, using tweet-retweet propagation network, in which all malicious profiles/bot will be eliminated and then spread and propagation scores will be calculated and saved as tweet

feature (see feature no: 46 and 47 in table 15). The spread and propagation pattern of the message is an important indicator of credibility assessment. After calculating all user-level features and tweet-level features, different machine learning models could be executed over these features for the prediction of post label/ rank / score. Therefore our model is following a hybrid approach combining both graph-based methods and feature-based methods.

**12.3. Post Credibility Score:** It could easily be observed that our proposed list of features completely covers all quality-related aspects of a tweet. These quality-related aspects

were mostly missed from the majority of studies in the literature. After calculating all user-level features and tweet-level features (see the recommended list of features in table 15) either any conventional Machine Learning Regression model (e.g.: Gradient Boosting, Ada-Boost, CAT Boost, LightGBM, SVM, Random Forest, Linear Regression, etc.) or any modern Learn to Rank model (e.g.: Lambda Rank, SVM Rank, Lambda Mart, etc.) could be executed to predict tweet's credibility score.

**12.4. User Level Scores:** Referring to our recommended solution, in addition to the basic tweet credibility score, different user level scores could easily be calculated based on the tweets of the user. It has been discussed earlier that many post-level features could compute user-level features. Examples of such User level scores are as follows. Computation of all such scores at the user level will implicitly reduce the dissemination of low-credibility contents, over microblogs.

**1. User %age of fake produced and propagated:** which will be a historic feature computed through no of fake tweets produced or propagated by that user.

**2. User Avg. Spread and Propagation Scores:** which will also be historic features, computed through avg. of all tweets spread and propagation scores of the user.

**3. User Avg. Credibility Score:** similarly User Avg. Credibility Score will be calculated by taking avg. of credibility score of all tweets of that user.

**4. User Top 3 Domains:** it could also be computed through all tweet topics tweeted by that user and the top 3 could be accumulated.

**12.5. Scores Convergence:** Above all user-level, scores and post-level scores could easily be calculated in real-time and displayed at respective entity levels. The chain of narrators is extremely important in assessing the message's credibility. Once a post is identified as fake then its producer must be penalized by incrementing its fake producer counter. Similarly, each fake propagator involved in post propagation within the post's chain of narrators must also be updated. It is worth mentioning that every post's credibility score will affect the respective user-level score and user score will also be affected by its post credibility. For example tweet's final credibility score will only be accumulated through all its chain of narrators and vice versa.

### XIII. FUTURE RESEARCH DIRECTIONS:

There is a need for benchmark/gold-standard credibility dataset construction. The dataset will include different forms of deceptions [33], like rumor, fake news, spam & scam, hoax, click-bait, junk science, conspiracy, and different forms of smear campaigns, etc. The dataset must also be enriched with hate speech, with its related concepts like abusive language, offensive language, general hate, cyberbullying, discrimination, flaming, harassment, profanity, toxic language or comment, extremism, radicalization, etc. [14]. There should also be sufficient malicious profiles (e.g. Bots/Cyborgs, etc.).

It must contain a good mix of news and non-news pieces of information. The dataset tagging should be done exactly in the way which is presented in section XII's heading 'Guided Data Tagging'. Regarding the features of the dataset. The following necessary features must be included for credibility assessment. The dataset must have a three-degree friends network (followers/following directed graph), user profiles, complete tweets of all users involved in the datasets with the number of replies & number favorites & who has favorited, etc., in addition to actual tweets which will be considered for credibility assessment. Actual and complete tweet-retweet multi-level propagation network (generally Twitter API provides flat retweeter's list), information of the list/ groups, media files, etc. The dataset's post should have a balanced number of domains e.g.: Politics, Entertainment, Sports, Education, etc. The dataset should also be developed through multiple microblogs and in different languages.

There are many challenges involved in the development of such a dataset because of accessibility privileges, the huge amount of data collection and management, strict tagging requirements, etc. Fortunately, there are few components of such dataset that are already available (see section VII) that need to be compiled concerning credibility, and missing components will be added.

In addition to the real-world labeled dataset. We need to implement the recommended system presented in this study, for its efficacy and performance evaluation.

After the necessary understanding of information credibility for microblogs presented through this study. There is a need to explore the literature regarding information credibility using multi-modal data and, explainable credibility assessment methods. It is very important that whatever credibility assessment is done by the system needs to be explained, that how the contents are categorized as not-credible or credible. Similarly, credibility assessment should make use of voice, image, and video from the post, in addition to text. Regarding the challenges and limitations, which are presented in different sections of the study therefore not discussed separately.

### XIV. CONCLUSION:

An effort of presenting the anatomy of information credibility for social media and microblogs was made, through a detailed and, organized study. Many research studies were conducted to assess automatic microblog's credibility but the majority of them had different concepts of credibility. Credibility is multi-disciplinary, hence there was no generalized or accepted credibility concept with all its necessary and detailed constructs/components. Therefore, it was necessary to understand the complete concept of information credibility from different disciplines. It could be accomplished through an organized study of all the problem dimensions and identification of comprehensive and necessary credibility constructs under credibility's definition. Such literature exploration and the fundamental study was missing regard-

ing the work done. Therefore to consolidate, standardized, identify gaps, propose solutions and recommendations in this area. We deeply explore the existing literature first, categories them along various dimensions, identify gaps and shortcomings then suggest important recommendations. As a result of a successful explorational study, a complete information credibility framework for social media is proposed. It is the first framework considering all necessary constructs of credibility identified in this study. Afterward, the presented framework is also transformed for microblogs credibility assessment. The transformation is done to individual features level for understanding and clarity. Therefore the framework can simply be implemented as a successful system. Another important aspect which we noticed missing in previous researches and therefore proposed, that Credibility should be measured through the narrators and narrations both considering their important aspects or bases of assessments. The narrator's assessment should be done on multiple bases such as its genuine social network influence, should always be truthful and unbiased, its area of expertise, popularity, and good reputation, etc. Similarly, narration could be assessed on its quality basis like it must be true, clean from spam & scams, rumors, and smear campaigns, etc. It should be informative, clear from the variety of hate speeches, and extreme biases, etc. Our credibility framework is based on both user and post. Which could provide two-fold benefits information credibility ratings as well as user credibility ratings. Later credibility (user credibility ratings) will be extremely helpful in other applications for example to assess the reviews of credible authors, considerations of credible user's recommendations, etc.

## Appendix: Terms Defined

**Claim:** Un-verified piece of news/ article/ information/ opinion in question, which could be rumor, hoax, satire, and fake news, etc.

**Fact-Checking:** Process of claim evaluation through authentic publish media, journalists, and domain experts, etc., and resulted as Fake, Real, etc.

**Satire:** is characterized by humor, irony, absurdity, exaggeration, and ridicule. They can mimic genuine news, primarily written to criticize.

**Hoax:** Deliberately fabricated falsehood made to masquerade as truth, intentionally conceived to deceive readers.

**Propaganda:** Information that tries to influence the emotion of the opinions, and the actions of target audiences through deceptive, selectively omitting, and one-sided messages. The purpose could be political, ideological, or religious, etc.

**Rumor:** Claim that has not been verified (may be true or false), apparently credible but hard to verify and spread from one person to another.

**Click-bait:** Low-quality journalism intended to attract traffic and monetize via advertising revenue.

**Meme:** A piece of information that replicates among people (Dawkins 1989). It bears similarities to infectious diseases

as both travel through social ties from one person to another. Piece of information mostly spread widely on the internet, often altered for humorous effect. Meme types are hashtags, URLs, Mentions, and Phrases.

**Astroturfing:** A particular type of abuse disguised as spontaneous "grassroots" behavior, but that is in reality carried out by a single person or organization. Non-genuine public support of an issue. Quiet related to spam.

**Sybil's:** Suspicious accounts, no malicious contents are posted, creating many fake identities to unfairly increase the power or influence of someone, therefore, produce a false sense of credibility. This concept is called Link Farming. Some similar terms to Sybil's are also popular e.g.: Sock-puppet, Zombie Followers, and Fake followers, etc.

**Bots, Trolls, Cyborg:** "Bots" are fully automated accounts and completely distinct from professional "trolls", which are human-run accounts, and the "Cyborg" accounts which combine human-generated content with automated posting.

**Botnets:** connected bots network.

**Social Spambots:** More sophisticated bots, mimic human-like behavior.

**Spambots/Content Polluters:** Traditional and simple type of bots, e.g.: Duplicate Spammers, Malicious Promoters, Self-Promoters, Friend Infiltrator, etc.

**Coordinated Behavior:** Chain of users which are developed to perform some pre-defined task of their master (example of pre-defined task could be: always like the post, add specific hashtag and mention, then forward post to others).

**Followers Fallacy:** Users with manipulated followers count. These untrustworthy users use bot activities to increase followers count for having high influence, popularity, or reputation. There are different ways, like online black-market services, they help the users to increase their followers/likes. Users can purchase bulk followers and likes from these markets. Users exploit such services to inflate followers, likes, and shares of the post to become more influential and popular.

**Extreme Bias:** Piece of information come from a particular point of view and may rely on propaganda, decontextualized information, and opinions distorted as facts.

**Linguistic Inquiry and Word Count (LIWC):** Psycholinguistic features are very important in credibility analysis through text, which could be computed by LIWC. It is a text analysis lexicon and a program that calculates the percentage of words in a given text that fall into one or more of over 80 linguistic, psychological, and topical categories indicating various social, cognitive, and affective processes. i.e.: the word 'cried' is part of four-word categories: sadness, negative emotion, overall affect, and a past tense verb.

## REFERENCES

- [1] Elisa Shearer and BY Jeffrey Gottfried. News use across social media platforms 2017. pew research center (2017), 17, 2017.
- [2] Andrew Perrin. Social media usage. Pew research center, pages 52–68, 2015.
- [3] Marcelo Mendoza, Barbara Poblete, and Carlos Castillo. Twitter under



- crisis: Can we trust what we rt? In Proceedings of the first workshop on social media analytics, pages 71–79, 2010.
- [4] Erica J Briscoe, Darren Scott Appling, and Heather Hayes. Social network derived credibility. In Recommendation and Search in Social Networks, pages 59–75. Springer, 2015.
- [5] Akshay Java, Xiaodan Song, Tim Finin, and Belle Tseng. Why we twitter: understanding microblogging usage and communities. In Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis, pages 56–65, 2007.
- [6] A LENHART. Teens, privacy and online social networks. how teens manage their online identities in the age of myspace', pew internet net & american life project report 2007. [http://www.pewinternet.org/PPF/r/211/report\\_display.asp](http://www.pewinternet.org/PPF/r/211/report_display.asp), Accessed: Aug. 12, 2021.
- [7] Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. What is twitter, a social network or a news media? In Proceedings of the 19th international conference on World wide web, pages 591–600, 2010.
- [8] Alfred Hermida. Twittering the news: The emergence of ambient journalism. Journalism practice, 4(3):297–308, 2010.
- [9] Sam Laird. How social media is taking over the news industry. Mashable. Retrieved on, 2012.
- [10] Nic Newman, Richard Fletcher, Antonis Kalogeropoulos, David Levy, and Rasmus-Kleis Nielsen. Reuters institute digital news report 2017. reuters institute; university of oxford, 2017.
- [11] Daron Acemoglu, Asuman Ozdaglar, and Ali ParandehGheibi. Spread of (mis) information in social networks. Games and Economic Behavior, 70(2):194–227, 2010.
- [12] Priyanka Meel and Dinesh Kumar Vishwakarma. Fake news, rumor information pollution in social media and web: A contemporary survey of state-of-the-arts, challenges and opportunities. Expert Systems with Applications, 153:112986, 2020.
- [13] Botambu Collins, Dinh Tuyen Hoang, Ngoc Thanh Nguyen, and Dosam Hwang. Trends in combating fake news on social media—a survey. Journal of Information and Telecommunication, 5(2):247–266, 2021.
- [14] Paula Fortuna and Sérgio Nunes. A survey on automatic detection of hate speech in text. ACM Computing Surveys (CSUR), 51(4):1–30, 2018.
- [15] K. A. Qureshi and M. Sabih. Un-compromised credibility: Social media based multi-class hate speech classification for text. IEEE Access, doi:10.1109/ACCESS.2021.3101977, 2021.
- [16] Jacob Ratkiewicz, Michael Conover, Mark Meiss, Bruno Gonçalves, Snehal Patil, Alessandro Flammini, and Filippo Menczer. Detecting and tracking the spread of astroturf memes in microblog streams. arXiv preprint arXiv:1011.3768, 2010.
- [17] Kai Shu, Suhang Wang, and Huan Liu. Beyond news contents: The role of social context for fake news detection. In Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, pages 312–320, 2019.
- [18] Majd Latah. The art of social bots: A review and a refined taxonomy. arXiv preprint arXiv:1905.03240, 2019.
- [19] Meeyoung Cha, Hamed Haddadi, Fabricio Benevenuto, and Krishna P Gummadi. Measuring user influence in twitter: The million follower fallacy. In fourth international AAAI conference on weblogs and social media, 2010.
- [20] Chengcheng Shao, Pik-Mai Hui, Lei Wang, Xinwen Jiang, Alessandro Flammini, Filippo Menczer, and Giovanni Luca Ciampaglia. Anatomy of an online misinformation network. PloS one, 13(4):e0196087, 2018.
- [21] Majed AlRubaian, Muhammad Al-Qurishi, Mabrook Al-Rakhami, Sk Md Mizanur Rahman, and Atif Alamri. A multistage credibility analysis model for microblogs. In 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pages 1434–1440. IEEE, 2015.
- [22] Soo Young Rieh, Meredith Ringel Morris, Miriam J Metzger, Helen Francke, and Grace YoungJoo Jeon. Credibility perceptions of content contributors and consumers in social media. Proceedings of the American Society for Information Science and Technology, 51(1):1–4, 2014.
- [23] Onur Varol, Emilio Ferrara, Clayton A Davis, Filippo Menczer, and Alessandro Flammini. Online human-bot interactions: Detection, estimation, and characterization. In Eleventh international AAAI conference on web and social media, 2017.
- [24] Christian Grimme, Mike Preuss, Lena Adam, and Heike Trautmann. Social bots: Human-like by means of human control? Big data, 5(4):279–293, 2017.
- [25] Hunt Allcott and Matthew Gentzkow. Social media and fake news in the 2016 election. Journal of economic perspectives, 31(2):211–36, 2017.
- [26] Matthew Hindman and Vlad Barash. Disinformation, and influence campaigns on twitter. 2018.
- [27] Robert Faris, Hal Roberts, Bruce Etling, Nikki Bourassa, Ethan Zuckerman, and Yochai Benkler. Partisanship, propaganda, and disinformation: Online media and the 2016 us presidential election. Berkman Klein Center Research Publication, 6, 2017.
- [28] Alexandre Bovet and Hernán A Makse. Influence of fake news in twitter during the 2016 us presidential election. Nature communications, 10(1):1–14, 2019.
- [29] Marco Viviani and Gabriella Pasi. Credibility in social media: opinions, news, and health information—a survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 7(5):e1209, 2017.
- [30] John ODonovan, Byungkyu Kang, Greg Meyer, Tobias Höllerer, and Sibel Adalii. Credibility in context: An analysis of feature distributions in twitter. In 2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Conference on Social Computing, pages 293–301. IEEE, 2012.
- [31] Alper Gün and Pinar Karagöz. A hybrid approach for credibility detection in twitter. In International Conference on Hybrid Artificial Intelligence Systems, pages 515–526. Springer, 2014.
- [32] Shamanth Kumar, Fred Morstatter, and Huan Liu. Twitter data analytics. Springer, 2014.
- [33] Esma Aïmeur, Hicham Hage, and Sabrine Amri. The scourge of online deception in social networks. In 2018 International Conference on Computational Science and Computational Intelligence (CSCI), pages 1266–1271. IEEE, 2018.
- [34] Fabián Riquelme and Pablo González-Cantergiani. Measuring user influence on twitter: A survey. Information processing & management, 52(5):949–975, 2016.
- [35] W Sherchan, S Nepal, and C Paris. A survey of trust in social networks. ACM Computing Surveys (CSUR), 45(4):1–33, 2013.
- [36] Eleonora Ciceri, Roman Fedorov, Eric Umhoza, Marco Brambilla, and Piero Fraternali. Assessing online media content trustworthiness, relevance and influence: an introductory survey. In KDWeb, pages 29–40, 2015.
- [37] Jun An, Wei Jiang Li, Lei Na Ji, and Feng Wang. A survey on information credibility on twitter. In Applied Mechanics and Materials, volume 401, pages 1788–1791. Trans Tech Publ, 2013.
- [38] Majed AlRubaian, Muhammad Al-Qurishi, Atif Alamri, Mabrook Al-Rakhami, Mohammad Mehdi Hassan, and Giancarlo Fortino. Credibility in online social networks: A survey. IEEE Access, 7:2828–2855, 2018.
- [39] C. C. Self. Credibility." in m. b. salwen & d. w. stacks (eds.), an integrated approach to communication theory and research. mahwah, nj: Lawrence erlbaum. 1996.
- [40] Credibility definition oxford, Accessed: Aug. 12, 2021. <https://en.oxforddictionaries.com/definition/credibility>.
- [41] Credibility definition merriam webster, Accessed: Aug. 12, 2021. <https://www.merriam-webster.com/dictionary/credibility>.
- [42] Pamela J Kalbfleisch. Credibility for the 21st century: Integrating perspectives on source, message, and media credibility in the contemporary media environment. In Communication yearbook 27, pages 307–350. Routledge, 2003.
- [43] Brian J Fogg. Prominence-interpretation theory: Explaining how people assess credibility online. In CHI'03 extended abstracts on human factors in computing systems, pages 722–723, 2003.
- [44] Wonchan Choi and Besiki Stvilia. Web credibility assessment: Conceptualization, operationalization, variability, and models. Journal of the Association for Information Science and Technology, 66(12):2399–2414, 2015.
- [45] Soo Young Rieh. Credibility and cognitive authority of information. 2010.
- [46] Brian J Fogg. Persuasive technology: using computers to change what we think and do. Ubiquity, 2002(December):2, 2002.
- [47] Shawn Tseng and BJ Fogg. Credibility and computing technology. Communications of the ACM, 42(5):39–44, 1999.
- [48] Brian Cugelman, Mike Thelwall, and Phil Dawes. Website credibility, active trust and behavioural intent. In International Conference on Persuasive Technology, pages 47–57. Springer, 2008.
- [49] Robin L Wakefield and Dwayne Whitten. Examining user perceptions of third-party organizations credibility and trust in an e-retailer. Journal of Organizational and End User Computing (JOEUC), 18(2):1–19, 2006.

- [50] David C Arnott, David Wilson, and Christina Sichtmann. An analysis of antecedents and consequences of trust in a corporate brand. *European Journal of Marketing*, 2007.
- [51] Gilles Sahut and André Tricot. Wikipedia: an opportunity to rethink the links between sources' credibility, trust and authority. *First Monday* 22(11), 2017.
- [52] Burns W Roper. *Public Attitudes Toward Television and Other Media in a Time of Change: The 14. Report in a Series by the Roper Organization Inc.* Television Information Office, 1985.
- [53] Miriam J Metzger, Andrew J Flanagin, Keren Eyal, Daisy R Lemus and Robert M McCann. Credibility for the 21st century: Integrating perspectives on source, message, and media credibility in the contemporary media environment. *Annals of the International Communication Association*, 27(1):293–335, 2003.
- [54] Thomas J Johnson and Barbara K Kaye. Cruising is believing?: Comparing internet and traditional sources on media credibility measures. *Journalism & Mass Communication Quarterly*, 75(2):325–340, 1998.
- [55] John Mashek, Larry McGill, and Adam C Powell. Lethargy '96: How the media covered a listless campaign. *Freedom Forum First Amendment Center at Vanderbilt University*, 1997.
- [56] Patrick Wilson. *Second-hand knowledge: An inquiry into cognitive authority*. 1983.
- [57] Soojung Kim. Questioners' credibility judgments of answers in a social question and answer site. *Information Research*, 15(2):15–2, 2010.
- [58] Soo Young Rieh. Judgment of information quality and cognitive authority in the web. *Journal of the American society for information science and technology*, 53(2):145–161, 2002.
- [59] C Nadine Wathen and Jacquelyn Burkell. Believe it or not: Factors influencing credibility on the web. *Journal of the American society for information science and technology*, 53(2):134–144, 2002.
- [60] S Shyam Sundar. Technology and credibility: Cognitive heuristics cued by modality, agency, interactivity and navigability. *Digital media, youth and learning*, pages 73–100, 2007.
- [61] Richard E Petty and John T Cacioppo. The elaboration likelihood model of persuasion. In *Communication and persuasion*, pages 1–24. Springer, 1986.
- [62] Shelly Chaiken. The heuristic model of persuasion. In *Social influence: the ontario symposium, volume 5*, pages 3–39, 1987.
- [63] Richard M Shiffrin and Walter Schneider. Controlled and automatic human information processing: II. perceptual learning, automatic attending and a general theory. *Psychological review*, 84(2):127, 1977.
- [64] Joseph B Walther. Interpersonal effects in computer-mediated interaction: A relational perspective. *Communication research*, 19(1):52–90, 1992.
- [65] Joseph B Walther. Social information processing theory. *Engaging theories in interpersonal communication: Multiple perspectives*, 391, 2008.
- [66] Miriam J Metzger. Making sense of credibility on the web: Models for evaluating online information and recommendations for future research. *Journal of the American society for information science and technology* 58(13):2078–2091, 2007.
- [67] Brian Hilligoss and Soo Young Rieh. Developing a unifying framework of credibility assessment: Construct, heuristics, and interaction in context. *Information Processing & Management*, 44(4):1467–1484, 2008.
- [68] Soo Young Rieh, Yong-Mi Kim, Ji Yeon Yang, and Beth St. Jean. A diary study of credibility assessment in everyday life information activities on the web: Preliminary findings. *Proceedings of the American Society for Information Science and Technology*, 47(1):1–10, 2010.
- [69] Adrien Friggeri, Lada Adamic, Dean Eckles, and Justin Cheng. Rumor cascades. In *Eighth International AAAI Conference on Weblogs and Social Media*, 2014.
- [70] Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. Prominent features of rumor propagation in online social media. In *2013 IEEE 13th International Conference on Data Mining*, pages 1103–1108. IEEE, 2013.
- [71] Soroush Vosoughi, Mostafa 'Neo' Mohsenvand, and Deb Roy. Rumor gauge: Predicting the veracity of rumors on twitter. *ACM transactions on knowledge discovery from data (TKDD)*, 11(4):1–36, 2017.
- [72] Eunsoo Seo, Prasant Mohapatra, and Tarek Abdelzaher. Identifying rumors and their sources in social networks. In *Ground/air multisensor interoperability, integration, and networking for persistent ISR III*, volume 8389, page 8389I1. International Society for Optics and Photonics, 2012.
- [73] Vahed Qazvinian, Emily Rosengren, Dragomir R Radev, and Qiaozhu Mei. Rumor has it: Identifying misinformation in microblogs. In *Proceedings of the conference on empirical methods in natural language processing*, pages 1589–1599. Association for Computational Linguistics, 2011.
- [74] Soroush Vosoughi. *Automatic detection and verification of rumors on Twitter*. PhD thesis, Massachusetts Institute of Technology, 2015.
- [75] Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J Jansen, Kam-Fai Wong, and Meeyoung Cha. Detecting rumors from microblogs with recurrent neural networks. 2016.
- [76] Arkaitz Zubiaga, Maria Liakata, Rob Procter, Geraldine Wong Sak Hoi, and Peter Tolmie. Analysing how people orient to and spread rumors in social media by looking at conversational threads. *PloS one*, 11(3), 2016.
- [77] Laura Sydell. We tracked down a fake-news creator in the suburbs. here's what we learned. *National Public Radio*, 23, 2016.
- [78] Chengcheng Shao, Giovanni Luca Ciampaglia, Onur Varol, Kai-Cheng Yang, Alessandro Flammini, and Filippo Menczer. The spread of low-credibility content by social bots. *Nature communications*, 9(1):1–9, 2018.
- [79] Adam Kucharski. Study epidemiology of fake news. *Nature*, 540(7634):525–525, 2016.
- [80] Erdem Beğenilmiş and Suzan Uskudarli. Organized behavior classification of tweet sets using supervised learning methods. In *Proceedings of the 8th International Conference on Web Intelligence, Mining and Semantics*, pages 1–9, 2018.
- [81] Soroush Vosoughi, Deb Roy, and Sinan Aral. The spread of true and false news online. *Science*, 359(6380):1146–1151, 2018.
- [82] Soroush Vosoughi, Mostafa 'Neo' Mohsenvand, and Deb Roy. Rumor gauge: Predicting the veracity of rumors on twitter. *ACM transactions on knowledge discovery from data (TKDD)*, 11(4):1–36, 2017.
- [83] Srijan Kumar and Neil Shah. False information on web and social media: A survey. *arXiv preprint arXiv:1804.08559*, 2018.
- [84] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context and spatiotemporal information for studying fake news on social media. *arXiv preprint arXiv:1809.01286*, 2018.
- [85] Bilal Ghanem, Paolo Rosso, and Francisco Rangel. Stance detection in fake news a combined feature representation. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 66–71, 2018.
- [86] Jacob Ratkiewicz, Michael Conover, Mark Meiss, Bruno Gonçalves, Snehal Patil, Alessandro Flammini, and Filippo Menczer. Truthy: mapping the spread of astroturf in microblog streams. In *Proceedings of the 20th international conference companion on World wide web*, pages 249–252, 2011.
- [87] Jacob Ratkiewicz, Michael D Conover, Mark Meiss, Bruno Gonçalves, Alessandro Flammini, and Filippo Menczer. Detecting and tracking political abuse in social media. In *Fifth international AAAI conference on weblogs and social media*, 2011.
- [88] Fabricio Benevenuto, Gabriel Magno, Tiago Rodrigues, and Virgilio Almeida. Detecting spammers on twitter. In *Collaboration, electronic messaging, anti-abuse and spam conference (CEAS)*, volume 6, page 12, 2010.
- [89] Sidharth Chhabra, Anupama Aggarwal, Fabricio Benevenuto, and Ponnuram Kumaraguru. Phi. social: the phishing landscape through short urls. In *Proceedings of the 8th Annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference*, pages 92–101, 2011.
- [90] Chris Grier, Kurt Thomas, Vern Paxson, and Michael Zhang. @ spam: the underground on 140 characters or less. In *Proceedings of the 17th ACM conference on Computer and communications security*, pages 27–37, 2010.
- [91] Sarita Yardi, Daniel Romero, Grant Schoenebeck, et al. Detecting spam in a twitter network. *First Monday*, 15(1), 2010.
- [92] Saptarshi Ghosh, Naveen Sharma, Fabricio Benevenuto, Niloy Ganguly, and Krishna Gummadi. Cognos: crowdsourcing search for topic experts in microblogs. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, pages 575–590, 2012.
- [93] Reyhan Yeniterzi and Jamie Callan. Constructing effective and efficient topic-specific authority networks for expert finding in social media. In *Proceedings of the first international workshop on Social media retrieval and analysis*, pages 45–50, 2014.
- [94] Sibel Adali, Fred Sisenda, and Malik Magdon-Ismael. Actions speak as loud as words: Predicting relationships from social behavior data. In

- 2331 Proceedings of the 21st international conference on World Wide Web<sup>2405</sup>  
 2332 pages 689–698, 2012. <sup>2406</sup>
- [95] Meng Jiang, Peng Cui, and Christos Faloutsos. Suspicious behavior<sup>2407</sup>  
 2333 detection: Current trends and future directions. *IEEE Intelligent Systems*<sup>2408</sup>  
 2334 31(1):31–39, 2016. <sup>2409</sup>
- [96] Stefano Cresci, Roberto Di Pietro, Marinella Petrocchi, Angelo Spog<sup>2410</sup>  
 2337 nardi, and Maurizio Tesconi. Social fingerprinting: detection of spambot<sup>2411</sup>  
 2338 groups through dna-inspired behavioral modeling. *IEEE Transactions on*<sup>2412</sup>  
 2339 *Dependable and Secure Computing*, 15(4):561–576, 2017. <sup>2413</sup>
- [97] Meng Jiang, Peng Cui, Alex Beutel, Christos Faloutsos, and Shiqian<sup>2414</sup>  
 2341 Yang. Catching synchronized behaviors in large networks: A graph<sup>2415</sup>  
 2342 mining approach. *ACM Transactions on Knowledge Discovery from*<sup>2416</sup>  
 2343 *Data (TKDD)*, 10(4):1–27, 2016. <sup>2417</sup>
- [98] Zi Chu, Steven Gianvecchio, Haining Wang, and Sushil Jajodia. De<sup>2418</sup>  
 2345 tecting automation of twitter accounts: Are you a human, bot, or cyborg<sup>2419</sup>  
 2346 *IEEE Transactions on Dependable and Secure Computing*, 9(6):811–824<sup>2420</sup>  
 2347 2012. <sup>2421</sup>
- [99] Chengcheng Shao, Giovanni Luca Ciampaglia, Onur Varol, Kai-Chen<sup>2422</sup>  
 2349 Yang, Alessandro Flammini, and Filippo Menczer. The spread of low<sup>2423</sup>  
 2350 credibility content by social bots. *Nature communications*, 9(1):1–9<sup>2424</sup>  
 2351 2018. <sup>2425</sup>
- [100] Clayton Allen Davis, Onur Varol, Emilio Ferrara, Alessandro Flammini,<sup>2426</sup>  
 2353 and Filippo Menczer. Botnot: A system to evaluate social bots. In<sup>2427</sup>  
 2354 *Proceedings of the 25th international conference companion on world*<sup>2428</sup>  
 2355 *wide web*, pages 273–274, 2016. <sup>2429</sup>
- [101] Kyumin Lee, Brian David Eoff, and James Caverlee. Seven months with<sup>2430</sup>  
 2357 the devils: A long-term study of content polluters on twitter. In *Fifth*<sup>2431</sup>  
 2358 *international AAAI conference on weblogs and social media*, 2011. <sup>2432</sup>
- [102] Leo G Stewart, Ahmer Arif, and Kate Starbird. Examining trolls and<sup>2433</sup>  
 2360 polarization with a retweet network. In *Proc. ACM WSDM, workshop*<sup>2434</sup>  
 2361 *on misinformation and misbehavior mining on the web*, 2018. <sup>2435</sup>
- [103] Patxi Galán-García, José Gaviria de la Puerta, Carlos Laorden Gómez,<sup>2436</sup>  
 2363 Igor Santos, and Pablo García Bringas. Supervised machine learning for<sup>2437</sup>  
 2364 the detection of troll profiles in twitter social network: Application to a<sup>2438</sup>  
 2365 real case of cyberbullying. *Logic Journal of the IGPL*, 24(1):42–53, 2016.<sup>2439</sup>
- [104] Eytan Bakshy, Jake M Hofman, Winter A Mason, and Duncan J Watts.<sup>2440</sup>  
 2367 Everyone’s an influencer: quantifying influence on twitter. In *Proceed*<sup>2441</sup>  
 2368 *ings of the fourth ACM international conference on Web search and data*<sup>2442</sup>  
 2369 *mining*, pages 65–74, 2011. <sup>2443</sup>
- [105] Adrien Guille, Hakim Hacid, Cecile Favre, and Djamel A Zighed. In<sup>2444</sup>  
 2371 formation diffusion in online social networks: A survey. *ACM Sigmod*<sup>2445</sup>  
 2372 *Record*, 42(2):17–28, 2013. <sup>2446</sup>
- [106] Mehrdad Farajtabar, Yichen Wang, Manuel Gomez-Rodriguez, Shuang<sup>2447</sup>  
 2374 Li, Hongyuan Zha, and Le Song. Coevolve: A joint point process model<sup>2448</sup>  
 2375 for information diffusion and network evolution. *The Journal of Machine*<sup>2449</sup>  
 2376 *Learning Research*, 18(1):1305–1353, 2017. <sup>2450</sup>
- [107] Pik-Mai Hui, Chengcheng Shao, Alessandro Flammini, Filippo Menczer,<sup>2451</sup>  
 2378 and Giovanni Luca Ciampaglia. The hoaxy misinformation and fact<sup>2452</sup>  
 2379 checking diffusion network. In *Twelfth International AAAI Conference*<sup>2453</sup>  
 2380 *on Web and Social Media*, 2018. <sup>2454</sup>
- [108] Ramanathan Guha, Ravi Kumar, Prabhakar Raghavan, and Andrew<sup>2455</sup>  
 2382 Tomkins. Propagation of trust and distrust. In *Proceedings of the 13th*<sup>2456</sup>  
 2383 *international conference on World Wide Web*, pages 403–412, 2004. <sup>2457</sup>
- [109] Mohsen Jamali and Martin Ester. Trustwalker: a random walk model for<sup>2458</sup>  
 2384 combining trust-based and item-based recommendation. In *Proceedings*<sup>2459</sup>  
 2385 *of the 15th ACM SIGKDD international conference on Knowledge*<sup>2460</sup>  
 2386 *discovery and data mining*, pages 397–406, 2009. <sup>2461</sup>
- [110] Wei Feng and Jianyong Wang. Retweet or not? personalized tweet re<sup>2462</sup>  
 2389 ranking. In *Proceedings of the sixth ACM international conference on*<sup>2463</sup>  
 2390 *Web search and data mining*, pages 577–586, 2013. <sup>2464</sup>
- [111] Srijith Ravikumar, Raju Balakrishnan, and Subbarao Kambhampati<sup>2465</sup>  
 2391 Ranking tweets considering trust and relevance. In *Proceedings of the*<sup>2466</sup>  
 2392 *Ninth International Workshop on Information Integration on the Web*<sup>2467</sup>  
 2393 *pages 1–4*, 2012. <sup>2468</sup>
- [112] Yajuan Duan, Long Jiang, Tao Qin, Ming Zhou, and Heung-Yeung Shum<sup>2469</sup>  
 2396 An empirical study on learning to rank of tweets. In *Proceedings of the*<sup>2470</sup>  
 2397 *23rd International Conference on Computational Linguistics*, pages 295–2470  
 2398 303. Association for Computational Linguistics, 2010. <sup>2471</sup>
- [113] Hongzhao Huang, Arkaitz Zubiaga, Heng Ji, Hongbo Deng, Dong Wang<sup>2472</sup>  
 2400 Hieu Le, Tarek Abdelzaher, Jiawei Han, Alice Leung, John Hancock<sup>2473</sup>  
 2401 et al. Tweet ranking based on heterogeneous networks. In *Proceeding*<sup>2474</sup>  
 2402 *of COLING 2012*, pages 1239–1256, 2012. <sup>2475</sup>
- [114] Jun Ito, Jing Song, Hiroyuki Toda, Yoshimasa Koike, and Satoshi Oyama<sup>2476</sup>  
 2403 Assessment of tweet credibility with lda features. In *Proceedings of*<sup>2477</sup>  
 2404 *the 24th International Conference on World Wide Web*, pages 953–958,  
 2015.
- [115] Shervin Malmasi and Marcos Zampieri. Detecting hate speech in social  
 media. arXiv preprint arXiv:1712.06427, 2017.
- [116] Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal,  
 Noura Farra, and Ritesh Kumar. Predicting the type and target of  
 offensive posts in social media. arXiv preprint arXiv:1902.09666, 2019.
- [117] Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek Bevendorff, and  
 Benno Stein. A stylometric inquiry into hyperpartisan and fake news.  
 arXiv preprint arXiv:1702.05638, 2017.
- [118] Amitabha Dey, Rafsan Zani Rafi, Shahriar Hasan Parash, Sauvik Kundu  
 Arko, and Amitabha Chakrabarty. Fake news pattern recognition using  
 linguistic analysis. In 2018 Joint 7th International Conference on  
 Informatics, Electronics & Vision (ICIEV) and 2018 2nd International  
 Conference on Imaging, Vision & Pattern Recognition (icIVPR), pages  
 305–309. IEEE, 2018.
- [119] Meredith Ringel Morris, Scott Counts, Asta Roseway, Aaron Hoff, and  
 Julia Schwarz. Tweeting is believing? understanding microblog credibil-  
 ity perceptions. In *Proceedings of the ACM 2012 conference on computer*  
*supported cooperative work*, pages 441–450, 2012.
- [120] Shafiza Mohd Shariff, Xiuzhen Zhang, and Mark Sanderson. On the  
 credibility perception of news on twitter: Readers, topics and features.  
*Computers in Human Behavior*, 75:785–796, 2017.
- [121] David Westerman, Patric R Spence, and Brandon Van Der Heide. A  
 social network as information: The effect of system generated reports of  
 connectedness on credibility on twitter. *Computers in Human Behavior*,  
 28(1):199–206, 2012.
- [122] Jiang Yang, Scott Counts, Meredith Ringel Morris, and Aaron Hoff.  
 Microblog credibility perceptions: comparing the usa and china. In  
*Proceedings of the 2013 conference on Computer supported cooperative*  
*work*, pages 575–586, 2013.
- [123] Sujoy Sikdar, Byungkyu Kang, John ODonovan, Tobias Höllerer, and  
 Sibel Adah. Understanding information credibility on twitter. In 2013  
 International Conference on Social Computing, pages 19–24. IEEE,  
 2013.
- [124] Carl I Hovland and Walter Weiss. The influence of source credibility on  
 communication effectiveness. *Public opinion quarterly*, 15(4):635–650,  
 1951.
- [125] Geoffrey Barbier and Huan Liu. Information provenance in social media.  
 In *International Conference on Social Computing, Behavioral-Cultural*  
*Modeling, and Prediction*, pages 276–283. Springer, 2011.
- [126] Kevin R Canini, Bongwon Suh, and Peter L Pirolli. Finding credible  
 information sources in social networks based on content and social  
 structure. In 2011 IEEE Third International Conference on Privacy,  
 Security, Risk and Trust and 2011 IEEE Third International Conference  
 on Social Computing, pages 1–8. IEEE, 2011.
- [127] Mohammad-Ali Abbasi and Huan Liu. Measuring user credibility  
 in social media. In *International Conference on Social Computing,*  
*Behavioral-Cultural Modeling, and Prediction*, pages 441–448. Springer,  
 2013.
- [128] Yuto Yamaguchi, Tsubasa Takahashi, Toshiyuki Amagasa, and Hiroyuki  
 Kitagawa. Turank: Twitter user ranking based on user-tweet graph analy-  
 sis. In *International conference on Web information systems engineering*,  
 pages 240–253. Springer, 2010.
- [129] David Westerman, Patric R Spence, and Brandon Van Der Heide. Social  
 media as information source: Recency of updates and credibility of  
 information. *Journal of computer-mediated communication*, 19(2):171–  
 183, 2014.
- [130] Lijuan Huang and Yeming Xiong. Evaluation of microblog users’  
 influence based on pagerank and users behavior analysis. 2013.
- [131] Byungkyu Kang, John O’Donovan, and Tobias Höllerer. Modeling  
 topic specific credibility on twitter. In *Proceedings of the 2012 ACM*  
*international conference on Intelligent User Interfaces*, pages 179–188,  
 2012.
- [132] Aditi Gupta and Ponnurangam Kumaraguru. Credibility ranking of  
 tweets during high impact events. In *Proceedings of the 1st workshop*  
*on privacy and security in online social media*, pages 2–8, 2012.
- [133] Hend S Al-Khalifa and Rasha M Al-Eidan. An experimental system for  
 measuring the credibility of news content in twitter. *International Journal*  
*of Web Information Systems*, 2011.
- [134] Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. Information  
 credibility on twitter. In *Proceedings of the 20th international conference*  
*on World wide web*, pages 675–684, 2011.


- [135] Xin Xia, Xiaohu Yang, Chao Wu, Shanping Li, and Linfeng Bao. Information credibility on twitter in emergency situation. In *Pacific Asia Workshop on Intelligence and Security Informatics*, pages 45–59. Springer, 2012.
- [136] Aditi Gupta, Ponnurangam Kumaraguru, Carlos Castillo, and Patrick Meier. Tweetcred: Real-time credibility assessment of content on twitter. In *International Conference on Social Informatics*, pages 228–243. Springer, 2014.
- [137] Krzysztof Lorek, Jacek Suehiro-Wiciński, Michal Jankowski-Lorek, and Amit Gupta. Automated credibility assessment on twitter. *Computer Science*, 16(2):157–168, 2015.
- [138] Majed Alrubaian, Muhammad Al-Qurishi, Mohammad Mehedi Hassan and Atif Alamri. A credibility analysis system for assessing information on twitter. *IEEE Transactions on Dependable and Secure Computing*, 15(4):661–674, 2016.
- [139] Jun Ito, Jing Song, Hiroyuki Toda, Yoshimasa Koike, and Satoshi Oyama. Assessment of tweet credibility with lda features. In *Proceedings of the 24th International Conference on World Wide Web*, pages 953–958. 2015.
- [140] Zhiwei Jin, Juan Cao, Yu-Gang Jiang, and Yongdong Zhang. News credibility evaluation on microblog with a hierarchical propagation model. In *2014 IEEE International Conference on Data Mining*, pages 230–239. IEEE, 2014.
- [141] Manish Gupta, Peixiang Zhao, and Jiawei Han. Evaluating event credibility on twitter. In *Proceedings of the 2012 SIAM International Conference on Data Mining*, pages 153–164. SIAM, 2012.
- [142] Jason RC Nurse, Ioannis Agraftiotis, Sadie Creese, Michael Goldsmith and Koen Lamberts. Building confidence in information-trustworthiness metrics for decision support. In *2013 12th IEEE International Conference on Trust, Security and Privacy in Computing and Communications*, pages 535–543. IEEE, 2013.
- [143] Amal Abdullah AlMansour and Costas S Iliopoulos. Using arabic microblogs features in determining credibility. In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 1212–1219, 2015.
- [144] Gillian Moran and Laurent Muzellec. ewom credibility on social networking sites: A framework. *Journal of Marketing Communications*, 23(2):149–161, 2017.
- [145] Jason RC Nurse, Syed Sadiqur Rahman, Sadie Creese, Michael Goldsmith, and Koen Lamberts. Information quality and trustworthiness: A topical state-of-the-art review. 2011.
- [146] Philip Meyer. Defining and measuring credibility of newspapers: Developing an index. *Journalism quarterly*, 65(3):567–574, 1988.
- [147] Frederick Fico, John D Richardson, and Steven M Edwards. Influence of story structure on perceived story bias and news organization credibility. *Mass communication & society*, 7(3):301–318, 2004.
- [148] Cecilie Gaziano and Kristin McGrath. Measuring the concept of credibility. *Journalism quarterly*, 63(3):451–462, 1986.
- [149] Miriam J Metzger, Andrew J Flanagin, Keren Eyal, Daisy R Lemus and Robert M McCann. Credibility for the 21st century: Integrating perspectives on source, message, and media credibility in the contemporary media environment. *Annals of the International Communication Association*, 27(1):293–335, 2003.
- [150] Keyi Xu, Ye Liu, Xueting Zhao, and Xiaoyue Dong. Trust them or not? a study on media credibility of newspapers accounts on sina weibo. *A Study on Media Credibility of Newspapers Accounts on Sina Weibo* (April 30, 2013), 2013.
- [151] Carl I Hovland and Walter Weiss. The influence of source credibility on communication effectiveness. *Public opinion quarterly*, 15(4):635–650, 1951.
- [152] JC McCroskey. *An introduction to communication in the classroom*. edina, mn: Burgess international group, 1992.
- [153] Jason J Teven and James C McCroskey. The relationship of perceived teacher caring with student learning and teacher evaluation. *Communication Education*, 46(1):1–9, 1997.
- [154] James C McCroskey and Jason J Teven. Goodwill: A reexamination of the construct and its measurement. *Communications Monographs*, 66(1):90–103, 1999.
- [155] DJ O’Keefe. *Theories of behavioral intention. Persuasion Theory & Research*, 2nd ed. Thousand Oaks, CA: Sage, pages 101–135, 2002.
- [156] Carl Iver Hovland, Irving Lester Janis, and Harold H Kelley. *Communication and persuasion*. 1953.
- [157] Alexandru L Ginsca, Adrian Popescu, Mihai Lupu, et al. Credibility in information retrieval. *Foundations and Trends® in Information Retrieval*, 9(5):355–475, 2015.
- [158] Yolanda Gil and Donovan Artz. Towards content trust of web resources. *Journal of Web Semantics*, 5(4):227–239, 2007.
- [159] Andrew J Flanagin and Miriam J Metzger. Digital media and youth: Unparalleled opportunity and unprecedented responsibility. *MacArthur Foundation Digital Media and Learning Initiative*, 2008.
- [160] Yu Suzuki. A credibility assessment for message streams on microblogs. In *2010 International Conference on P2P, Parallel, Grid, Cloud and Internet Computing*, pages 527–530. IEEE, 2010.
- [161] Byungkyu Kang, TH Höllerer, Matthew Turk, Xifeng Yan, and John O’Donovan. An analysis of credibility in microblogs. PhD thesis, MS thesis, Dept. Comput. Sci., Univ. California, Santa Barbara, Santa Barbara, 2012.
- [162] Mya Thandar and Sasiporn Usanavasin. Measuring opinion credibility in twitter. In *Recent Advances in Information and Communication Technology 2015*, pages 205–214. Springer, 2015.
- [163] Manish Gupta, Peixiang Zhao, and Jiawei Han. Evaluating event credibility on twitter. In *Proceedings of the 2012 SIAM International Conference on Data Mining*, pages 153–164. SIAM, 2012.
- [164] Tanushree Mitra and Eric Gilbert. Credbank: A large-scale social media corpus with associated credibility annotations. In *Ninth International AAAI Conference on Web and Social Media*, 2015.
- [165] Guy J Golan. *New perspectives on media credibility research*, 2010.
- [166] Soo Young Rieh and David R Danielson. *Credibility: A multidisciplinary framework*. 2007.
- [167] Richard Y Wang and Diane M Strong. Beyond accuracy: What data quality means to data consumers. *Journal of management information systems*, 12(4):5–33, 1996.
- [168] BJ Fogg, Jonathan Marshall, Othman Laraki, Alex Osipovich, Chris Varma, Nicholas Fang, Jyoti Paul, Akshay Rangnekar, John Shon, Preeti Swani, et al. What makes web sites credible? a report on a large quantitative study. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 61–68, 2001.
- [169] Traci Hong. The influence of structural and message features on web site credibility. *Journal of the American Society for Information Science and Technology*, 57(1):114–127, 2006.
- [170] David Robins and Jason Holmes. Aesthetics and credibility in web site design. *Information Processing & Management*, 44(1):386–399, 2008.
- [171] Soo Young Rieh and Brian Hilligoss. College students’ credibility judgments in the information-seeking process. *Digital media, youth, and credibility*, pages 49–72, 2008.
- [172] Shelly Chaiken and Yaacov Trope. *Dual-process theories in social psychology*. Guilford Press, 1999.
- [173] Francesco Pierri and Stefano Ceri. False news on social media: A data-driven survey. *ACM SIGMOD Record*, 48(2):18–27, 2019.
- [174] Jon Roozenbeek and Sander van der Linden. Fake news game confers psychological resistance against online misinformation. *Palgrave Communications*, 5(1):1–10, 2019.
- [175] Svitlana Volkova, Kyle Shaffer, Jin Yea Jang, and Nathan Hodas. Separating facts from fiction: Linguistic models to classify suspicious and trusted news posts on twitter. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 647–653, 2017.
- [176] Bin Guo, Yasan Ding, Lina Yao, Yunji Liang, and Zhiwen Yu. The future of false information detection on social media: New perspectives and trends. *ACM Computing Surveys (CSUR)*, 53(4):1–36, 2020.
- [177] Fang Jin, Wei Wang, Liang Zhao, Edward Dougherty, Yang Cao, Chang-Tien Lu, and Naren Ramakrishnan. Misinformation propagation in the age of twitter. *Computer*, (12):90–94, 2014.
- [178] Xiaoming Wang, Yaguang Lin, Yanxin Zhao, Lichen Zhang, Juhua Liang, and Zhipeng Cai. A novel approach for inhibiting misinformation propagation in human mobile opportunistic networks. *Peer-to-Peer Networking and Applications*, 10(2):377–394, 2017.
- [179] Natascha A Karlova and Karen E Fisher. A social diffusion model of misinformation and disinformation for understanding human information behaviour. 2013.
- [180] Xumeng Chen, Leo Yu-Ho Lo, and Huamin Qu. Sirenless: reveal the intention behind news. *arXiv preprint arXiv:2001.02731*, 2020.
- [181] Dilek Küçük and Fazli Can. Stance detection: A survey. *ACM Computing Surveys (CSUR)*, 53(1):1–37, 2020.

- [182] Eugenio Tacchini, Gabriele Ballarin, Marco L Della Vedova, Stefan Moret, and Luca de Alfaro. Some like it hoax: Automated fake news detection in social networks. arXiv preprint arXiv:1704.07506, 2017.
- [183] Credibility ground truths, 2020. <https://didyoucheckfirst.wordpress.com/opensource-co-news-sites/>.
- [184] Jennifer Golbeck, Cristina Robles, Michon Edmondson, and Karen Turner. Predicting personality from twitter. In 2011 IEEE third international conference on privacy, security, risk and trust and 2011 IEEE third international conference on social computing, pages 149–156. IEEE, 2011.
- [185] Daniele Quercia, Michal Kosinski, David Stillwell, and Jon Crowcroft. Our twitter profiles, our selves: Predicting personality with twitter. In 2011 IEEE third international conference on privacy, security, risk and trust and 2011 IEEE third international conference on social computing, pages 180–185. IEEE, 2011.
- [186] Michal Kakol, Radoslaw Nielek, and Adam Wierzbicki. Understanding and predicting web content credibility using the content credibility corpus. *Information Processing & Management*, 53(5):1043–1061, 2017.
- [187] Robert Thomson, Naoya Ito, Hinako Suda, Fangyu Lin, Yafei Liu, Ryo Hayasaka, Ryuzo Isochi, and Zian Wang. Trusting tweets: The fukushima disaster and information source credibility on twitter. In Proceedings of the 9th International ISCRAM Conference, pages 1–10. Vancouver: Simon Fraser University, 2012.
- [188] Marco Viviani and Gabriella Pasi. Quantifier guided aggregation for the veracity assessment of online reviews. *International Journal of Intelligent Systems*, 32(5):481–501, 2017.
- [189] Aditi Gupta, Hemank Lamba, Ponnuram Kumaraguru, and Anupam Joshi. Faking sandy: characterizing and identifying fake images on twitter during hurricane sandy. In Proceedings of the 22nd international conference on World Wide Web, pages 729–736, 2013.
- [190] Nurul H Idris, MJ Jackson, and MHI Ishak. A conceptual model of the automated credibility assessment of the volunteered geographic information formation. In IOP Conference Series: Earth and Environmental Science, volume 18, page 012070. IOP Publishing, 2014.
- [191] Rasha M BinSultan Al-Eidan, Henda S Al-Khalifa, and AbdulMalik S Al-Salman. Measuring the credibility of arabic text content in twitter. In 2010 Fifth International Conference on Digital Information Management (ICDIM), pages 285–291. IEEE, 2010.
- [192] Amal Abdullah AlMansour, Ljiljana Brankovic, and Costas S Iliopoulos. A model for recalibrating credibility in different contexts and languages: a twitter case study. *International Journal of Digital Information and Wireless Communications (IJDWC)*, 4(1):53–62, 2014.
- [193] Thomas J Johnson and Barbara K Kaye. Reasons to believe: Influence of credibility on motivations for using social networks. *Computers in human behavior*, 50:544–555, 2015.
- [194] Thomas J Johnson and Barbara K Kaye. Credibility of social network sites for political information among politically interested internet users. *Journal of Computer-mediated communication*, 19(4):957–974, 2014.
- [195] James Schaffer, Tarek Abdelzaher, Debra Jones, Tobias Höllerer, Cleotilde Gonzalez, Jason Harman, and John O'Donovan. Truth, lies and data: Credibility representation in data analysis. In 2014 IEEE International Inter-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support (CogSIMA), pages 28–34. IEEE, 2014.
- [196] Thomas J Johnson and Barbara K Kaye. The dark side of the boon: credibility, selective exposure and the proliferation of online sources of political information. *Computers in Human Behavior*, 29(4):1862–1871, 2013.
- [197] Q Vera Liao, Peter Pirolli, and Wai-Tat Fu. An act-r model of credibility judgment of micro-blogging web pages. *ICCM 2012 Proceedings*, 103:2754–2755, 2012.
- [198] Byungkyu Kang, Tobias Höllerer, and John O'Donovan. Believe it or not? analyzing information credibility in microblogs. In Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015, pages 611–616, 2015.
- [199] Miriam J Metzger and Andrew J Flanagin. Credibility and trust of information in online environments: The use of cognitive heuristics. *Journal of pragmatics*, 59:210–220, 2013.
- [200] Qin Gao, Ye Tian, and Mengyuan Tu. Exploring factors influencing chinese user's perceived credibility of health and safety information on weibo. *Computers in Human Behavior*, 45:21–31, 2015.
- [201] Chad Edwards, Autumn Edwards, Patric R Spence, and Ashleigh K Shelton. Is that a bot running the social media feed? testing the differences in perceptions of communication quality for a human agent and a bot agent on twitter. *Computers in Human Behavior*, 33:372–376, 2014.
- [202] Soo Young Rieh. Participatory web users' information activities and credibility assessment. 2010.
- [203] Erica J Briscoe, Darren Scott Appling, and Heather Hayes. Social network derived credibility. In *Recommendation and Search in Social Networks*, pages 59–75. Springer, 2015.
- [204] Tom Bakker, Damian Trilling, Claes de Vreese, Luzia Helfer, and Klaus Schönbach. The context of content: the impact of source and setting on the credibility of news. *Recherches en communication*, 40, 2013.
- [205] Kirsten A Johnson. The effect of twitter posts on students' perceptions of instructor credibility. *Learning, Media and Technology*, 36(1):21–38, 2011.
- [206] Robert Thomson, Naoya Ito, Hinako Suda, Fangyu Lin, Yafei Liu, Ryo Hayasaka, Ryuzo Isochi, and Zhou Wang. Trusting tweets: The fukushima disaster and information source credibility on twitter. In ISCRAM, 2012.
- [207] Arkaitz Zubiaga and Heng Ji. Tweet, but verify: epistemic study of information verification on twitter. *Social Network Analysis and Mining*, 4(1):163, 2014.
- [208] Jocelyn M DeGroot, Valerie J Young, and Sarah H VanSlette. Twitter use and its effects on student perception of instructor credibility. *Communication Education*, 64(4):419–437, 2015.
- [209] Jason RC Nurse, Ioannis Agraftotis, Michael Goldsmith, Sadie Creese, and Koen Lamberts. Two sides of the coin: measuring and communicating the trustworthiness of online information. *Journal of Trust Management*, 1(1):5, 2014.
- [210] Aditya Pal and Scott Counts. What's in a @ name? how name value biases judgment of microblog authors. In ICWSM. Citeseer, 2011.
- [211] Eun Go, Kyung Han You, Eunhwa Jung, and Hongjin Shim. Why do we use different types of websites and assign them different levels of credibility? structural relations among users' motives, types of websites, information credibility, and trust in the press. *Computers in Human Behavior*, 54:231–239, 2016.
- [212] Mi Rosie Jahng and Jeremy Littau. Interacting is believing: Interactivity, social cue, and perceptions of journalistic credibility on twitter. *Journalism & Mass Communication Quarterly*, 93(1):38–58, 2016.
- [213] Shafiza Mohd Shariff, Xiuzhen Zhang, and Mark Sanderson. User perception of information credibility of news on twitter. In European conference on information retrieval, pages 513–518. Springer, 2014.
- [214] Sujoy Kumar Sikdar, Byungkyu Kang, John O'Donovan, Tobias Hollerer, and Sibel Adal. Cutting through the noise: Defining ground truth in information credibility on twitter. *Human*, 2(3):151–167, 2013.
- [215] Suliman Aladhadh, Xiuzhen Zhang, and Mark Sanderson. Tweet author location impacts on tweet credibility. In proceedings of the 2014 Australasian document computing symposium, pages 73–76, 2014.
- [216] Eva Jaho, Efstratios Tzoannos, Aris Papadopoulos, and Nikos Sarris. Alethiometer: a framework for assessing trustworthiness and content validity in social media. In Proceedings of the 23rd International Conference on World Wide Web, pages 749–752, 2014.
- [217] Amal Abdullah AlMansour, Ljiljana Brankovic, and Costas S Iliopoulos. Evaluation of credibility assessment for microblogging: models and future directions. In Proceedings of the 14th International Conference on Knowledge Technologies and Data-driven Business, pages 1–4, 2014.
- [218] Rada Mihalcea and Dragomir Radev. Graph-based natural language processing and information retrieval. Cambridge university press, 2011.
- [219] William Yang Wang. "liar, liar pants on fire": A new benchmark dataset for fake news detection. arXiv preprint arXiv:1705.00648, 2017.
- [220] Kai Shu, Suhang Wang, and Huan Liu. Beyond news contents: The role of social context for fake news detection. In Proceedings of the twelfth ACM international conference on web search and data mining, pages 312–320, 2019.
- [221] Craig Silverman. This analysis shows how viral fake election news stories outperformed real news on facebook, buzzfeed news 2016. URL: <https://zenodo.org/record/1239675>, 16, Accessed: Aug. 12, 2021.
- [222] Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. Declare: Debunking fake news and false claims using evidence-aware deep learning. arXiv preprint arXiv:1809.06416, 2018.
- [223] Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra Lefevre, and Rada Mihalcea. Automatic detection of fake news. arXiv preprint arXiv:1708.07104, 2017.
- [224] Chengcheng Shao, Giovanni Luca Ciampaglia, Alessandro Flammini, and Filippo Menczer. Hoaxy: A platform for tracking online misinformation.

- 2768 mation. In Proceedings of the 25th international conference companion 2840  
 2769 on world wide web, pages 745–750, 2016. 2841
- 2770 [225] Megan Risdal. Fake news dataset 2017. URL: <https://www.kaggle.com/mrisdal/fake-news>, Accessed: Aug. 12, 2021. 2842
- 2771 2843
- 2772 [226] Leon Derczynski, Kalina Bontcheva, Maria Liakata, Rob Procter, Geraldine Wong Sak Hoi, and Arkaitz Zubiaga. Semeval-2017 task 8: Rumour evaluation: Determining rumour veracity and support for rumours. arXiv preprint arXiv:1704.05972, 2017. 2844
- 2773 2845
- 2774 2846
- 2775 2847
- 2776 [227] Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J Jansen, Kam-Fai Wong, and Meeyoung Cha. Detecting rumors from microblogs with recurrent neural networks. 2016. 2848
- 2777 2849
- 2778 2850
- 2779 [228] Francesco Pierri and Stefano Ceri. False news on social media: a data-driven survey. ACM Sigmod Record, 48(2):18–27, 2019. 2851
- 2780 2852
- 2781 [229] Alexandra Olteanu, Stanislav Peshterliev, Xin Liu, and Karl Aberer. Web credibility: Features exploration and credibility prediction. In European conference on information retrieval, pages 557–568. Springer, 2013. 2853
- 2782 2854
- 2783 2855
- 2784 [230] Steven A McCornack, Kelly Morrison, Jihyun Esther Paik, Amy M Wisner, and Xun Zhu. Information manipulation theory 2: A propositional theory of deceptive discourse production. Journal of Language and Social Psychology, 33(4):348–377, 2014. 2856
- 2785 2857
- 2786 2858
- 2787 2859
- 2788 [231] Oghenefejiro Winnie Makinde. Assessing the credibility of online social network messages. PhD thesis, University of Derby, 2018. 2860
- 2789 2861
- 2790 [232] Marcia K Johnson and Carol L Raye. Reality monitoring. Psychological review, 88(1):67, 1981. 2862
- 2791 2863
- 2792 [233] Miron Zuckerman, Bella M DePaulo, and Robert Rosenthal. Verbal and nonverbal communication of deception. In Advances in experimental social psychology, volume 14, pages 1–59. Elsevier, 1981. 2864
- 2793 2865
- 2794 2866
- 2795 [234] Udo Undeutsch. Beurteilung der glaubhaftigkeit von aussagen. Handbuch der psychologie, 11:26–181, 1967. 2867
- 2796 2868
- 2797 [235] Pamela M Homer and Lynn R Kahle. Source expertise, time of source identification, and involvement in persuasion: An elaborative processing perspective. Journal of advertising, 19(1):30–39, 1990. 2869
- 2798 2870
- 2799 2871
- 2800 [236] Ruohan Li and Ayoung Suh. Factors influencing information credibility on social media platforms: Evidence from facebook pages. Procedia computer science, 72:314–328, 2015. 2872
- 2801 2873
- 2802 2874
- 2803 [237] Blake E Ashforth and Fred Mael. Social identity theory and the organization. Academy of management review, 14(1):20–39, 1989. 2875
- 2804 2876
- 2805 [238] Gary Cronkrite and Jo Liska. A critique of factor analytic approaches to the study of credibility. Communications Monographs, 43(2):91–107, 1976. 2877
- 2806 2878
- 2807 2879
- 2808 [239] David Westerman, Patric R Spence, and Brandon Van Der Heide. Social media as information source: Recency of updates and credibility of information. Journal of computer-mediated communication, 19(2):171–183, 2014. 2880
- 2809 2881
- 2810 2882
- 2811 2883
- 2812 [240] Xiao Hu. Assessing Source Credibility On Social Media—An Electronic Word-Of-Mouth Communication Perspective. PhD thesis, Bowling Green State University, 2015. 2884
- 2813 2885
- 2814 2886
- 2815 [241] Kevin R Canini, Bongwon Suh, and Peter L Pirolli. Finding credible information sources in social networks based on content and social structure. In 2011 IEEE Third International Conference on Privacy, Security, Risk and Trust and 2011 IEEE Third International Conference on Social Computing, pages 1–8. IEEE, 2011. 2887
- 2816 2888
- 2817 2889
- 2818 2890
- 2819 2891
- 2820 [242] Johan Jessen and Anker Helms Jørgensen. Aggregated trustworthiness: Redefining online credibility through social validation. First Monday, 2012. 2892
- 2821 2893
- 2822 2894
- 2823 [243] Eleonora Ciceri, Roman Fedorov, Eric Umuhoza, Marco Brambilla, and Piero Fraternali. Assessing online media content trustworthiness, relevance and influence: an introductory survey. In KDWeb, pages 29–40, 2015. 2895
- 2824 2896
- 2825 2897
- 2826 2898
- 2827 [244] Damon Centola and Robb Wilier. The emperor’s dilemma. Theories of Social Order: A Reader, page 276, 2009. 2900
- 2828 2901
- 2829 2902
- 2830 [245] Morton Deutsch and Harold B Gerard. A study of normative and informational social influences upon individual judgment. The journal of abnormal and social psychology, 51(3):629, 1955. 2903
- 2831 2904
- 2832 2905
- 2833 [246] Timur Kuran and Cass R Sunstein. Availability cascades and risk regulation. Stan. L. Rev., 51:683, 1998. 2906
- 2834 2907
- 2835 [247] David Dunning, Dale W Griffin, James D Milojkovic, and Lee Ross. The overconfidence effect in social prediction. Journal of personality and social psychology, 58(4):568, 1990. 2908
- 2836 2909
- 2837 2910
- 2838 [248] Emily Pronin, Justin Kruger, Kenneth Savitsky, and Lee Ross. You don’t know me, but i know you: The illusion of asymmetric insight. Journal of Personality and Social Psychology, 81(4):639, 2001. 2911
- 2839 2912
- 2913
- [249] Andrew Ward, L Ross, E Reed, E Turiel, and T Brown. Naive realism in everyday life: Implications for social conflict and misunderstanding. Values and knowledge, pages 103–135, 1997.
- [250] Jonathan L Freedman and David O Sears. Selective exposure. In Advances in experimental social psychology, volume 2, pages 57–97. Elsevier, 1965.
- [251] Raymond S Nickerson. Confirmation bias: A ubiquitous phenomenon in many guises. Review of general psychology, 2(2):175–220, 1998.
- [252] Robert J Fisher. Social desirability bias and the validity of indirect questioning. Journal of consumer research, 20(2):303–315, 1993.
- [253] Harvey Leibenstein. Bandwagon, snob, and veblen effects in the theory of consumers’ demand. The quarterly journal of economics, 64(2):183–207, 1950.
- [254] Sai T Moturu and Huan Liu. Quantifying the trustworthiness of social media content. Distributed and Parallel Databases, 29(3):239–260, 2011.
- [255] Stephan Ten Kate. Trustworthiness within social networking sites: A study on the intersection of hci and sociology. Unpublished Business Studies Master, University of Amsterdam, Amsterdam, 2009.
- [256] Abdullah Algarni, Hashem AlMakrmi, and Abdulrahman Alarifi. Toward evaluating trustworthiness of social networking site users: Reputation-based method. Archives of Business Research, 7(3):27–41, 2019.
- [257] Lawrence E Boehm. The validity effect: A search for mediating variables. Personality and Social Psychology Bulletin, 20(3):285–293, 1994.
- [258] Péter Bálint and Géza Bálint. The semmelweis-reflex. Orvosi hetilap, 150(30):1430–1430, 2009.
- [259] Colin MacLeod, Andrew Mathews, and Philip Tata. Attentional bias in emotional disorders. Journal of abnormal psychology, 95(1):15, 1986.
- [260] Kathleen Hall Jamieson and Joseph N Cappella. Echo chamber: Rush Limbaugh and the conservative media establishment. Oxford University Press, 2008.
- [261] Carl I Hovland, OJ Harvey, and Muzafer Sherif. Assimilation and contrast effects in reactions to communication and attitude change. The Journal of Abnormal and Social Psychology, 55(2):244, 1957.
- [262] Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. In Handbook of the fundamentals of financial decision making: Part I, pages 99–127. World Scientific, 2013.
- [263] Nico H Frijda et al. The emotions. Cambridge University Press, 1986.
- [264] Julio Amador, Axel Oehmichen, and Miguel Molina-Solana. Characterizing political fake news in twitter by its meta-data. arXiv preprint arXiv:1712.05999, 2017.
- [265] Mansour Alsaleh, Abdulrahman Alarifi, Abdul Malik Al-Salman, Mohammed Alfayez, and Abdulmajeed Almuhsain. Tsd: Detecting sybil accounts in twitter. In 2014 13th International Conference on Machine Learning and Applications, pages 463–469. IEEE, 2014.
- [266] Wiki fake sites, 2020. [https://en.wikipedia.org/wiki/List\\_of\\_fake\\_news\\_websites](https://en.wikipedia.org/wiki/List_of_fake_news_websites).
- [267] News guard tech, 2020. <https://www.newsguardtech.com/>.
- [268] News media bias, 2020. <https://www.allsides.com/>.
- [269] Media bias fact check, 2020. <https://mediabiasfactcheck.com/>.
- [270] Stefano Cresci, Roberto Di Pietro, Marinella Petrocchi, Angelo Spognardi, and Maurizio Tesconi. The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race. In Proceedings of the 26th international conference on world wide web companion, pages 963–972, 2017.
- [271] Internettrusttoolnewsguardwebsite, Accessed: Aug. 12, 2021. <https://www.newsguardtech.com/>.
- [272] Yla R Tausczik and James W Pennebaker. The psychological meaning of words: Liwc and computerized text analysis methods. Journal of language and social psychology, 29(1):24–54, 2010.
- [273] Xinyi Zhou, Atishay Jain, Vir V Phoha, and Reza Zafarani. Fake news early detection: A theory-driven model. Digital Threats: Research and Practice, 1(2):1–25, 2020.
- [274] Marta Recasens, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky. Linguistic models for analyzing and detecting biased language. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1650–1659, 2013.
- [275] Factcheckingwebsiteslist, Accessed: Aug. 12, 2021. <https://reporterslab.org/fact-checking/>.
- [276] William H Dutton, Grant Blank, and Darja Groselj. Cultures of the internet: the internet in Britain: Oxford Internet Survey 2013 Report. Oxford Internet Institute, 2013.
- [277] Soroush Vosoughi, Deb Roy, and Sinan Aral. The spread of true and false news online. Science, 359(6380):1146–1151, 2018.


2914 [278] David MJ Lazer, Matthew A Baum, Yochai Benkler, Adam J Berinsky,  
2915 Kelly M Greenhill, Filippo Menczer, Miriam J Metzger, Brendan Nyhan,  
2916 Gordon Pennycook, David Rothschild, et al. The science of fake news.  
2917 Science, 359(6380):1094–1096, 2018.

2918  
2919  
2920  
2921  
2922  
2923  
2924  
2925  
2926  
2927  
2928  
2929




**KHUBAIB AHMED QURESHI** is currently affiliated with the Department of Computer Science, DHA Suffa University, Karachi, Pakistan, as Assistant Professor, and Head of Data Science Program. He was Head of Computer Science Department, HIMS, Hamdard University, Karachi, Pakistan. He is having 20 years of comprehensive research and teaching experience, continuing research in the area of Computer science named Data Science, Complex Networks, and Social Computing, etc. He has authored several research articles along with chapters in different books.

2930  
2931  
2932  
2933  
2934  
2935  
2936  
2937  
2938  
2939  
2940  
2941  
2942  
2943  
2944  
2945



**DR. RAUF AHMED SHAMS MALICK** received his PhD at University of Karachi. He has been visiting scholar at NIG(Japan), and UCLA (USA). He has founded several companies with state of the art products related to social media, location based analytics and organizational networks. He is currently involved in complex system research and pursuing problems in the area of biological networks, networked economics, and personality traits. He has distinguished background in designing novel solutions for complex systems. He is currently affiliated to Department of Computer Science, National University of Computer and Emerging Sciences, as Assistant Professor, continuing research in the specialized scientific area of Computer science, Complex Networks, Social Computing, Bioinformatics, Integrated system. He has authored several articles along with chapters in different books.

2946  
2947  
2948  
2949  
2950  
2951  
2952  
2953  
2954  
2955  
2956  
2957  
2958  
2959  
2960  
2961  
2962  
2963  
2964  
2965  
2966  
2967



**DR. MUHAMMAD SABIH** received the B.E. degree in industrial electronics from the IIEE-NED University, Pakistan, in 2000. He received his M.S. and PhD degrees in Systems Engineering from KFUPM, Dhahran, Saudi Arabia, in 2009 and 2014 respectively. During his research period at KFUPM, he worked for several applied research collaborations between KFUPM and MIT. His data driven research work expanded to a funded industry-academia collaboration project between KFUPM and Yokogawa-Saudi Arabia, and turned into a US patent. He also worked as Algorithm Specialist and developed anomaly detection algorithms using Python on the pipeline inspection data at leading Research and Technology Center (RTRC) of German ROSEN Group at Dhahran Techno-Valley (DTV), Saudi Arabia. He is working in the field of Computer and Electrical Engineering since 2009 and professional member of International Society of Automation (ISA). He is currently an Assistant Professor in DHA Suffa University and actively engaged in developing solutions from industrial data utilizing machine learning methods for estimation, modeling, and compensation. He has one US patent and around 10 peer-reviewed papers. His current research interest include Industry 4.0, Data Science, Modeling, and Estimation for real world problems.