

Language, Twitter and Academic Conferences

Ruth García¹ Diego Gómez² Denis Parra²
Christoph Trattner³ Andreas Kaltenbrunner¹ Eduardo Graells-Garrido⁴

¹Eurecat. Barcelona, Spain

²Pontificia Universidad Católica de Chile. Santiago, Chile

³NTNU. Trondheim, Norway

⁴Telefónica I+D. Santiago, Chile

ABSTRACT

Using Twitter during academic conferences is a way of engaging and connecting an audience inherently multicultural by the nature of scientific collaboration. English is expected to be the *lingua franca* bridging the communication and integration between native speakers of different mother tongues. However, little research has been done to support this assumption. In this paper we analyzed how integrated language communities are by analyzing the scholars' tweets used in 26 Computer Science conferences over a time span of five years. We found that although English is the most popular language used to tweet during conferences, a significant proportion of people also tweet in other languages. In addition, people who tweet solely in English interact mostly within the same group (English monolinguals), while people who speak other languages interact more with different *lingua groups*. Finally, we also found higher interaction between people tweeting in different languages. These results suggest a relation between the number of languages a user speaks and their interaction dynamics in online communities.

Categories and Subject Descriptors

J.4 [Social and Behavioral Sciences]: Sociology

Keywords

Twitter; culture; language; academic conferences

1. INTRODUCTION

In the past few years, Twitter has been used as a conference backchannel platform in academic events targeting the expansion of the community's communication and participation [1, 10]. Attendees using Twitter are generally involved in note taking, sharing resources and reporting individual real-time reactions to events, covering both conference presentations and conference social activities. This supports scholars' activities such as disseminating their work and engaging general public and newcomer scientists into the research communities [8]. It is a common practice in research conferences to use hashtags in the tweets to identify that particular event (e.g. #hypertext2015).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

HT'15 September 01 - 04, 2015, Guzelyurt, TRNC, Cyprus

©2015 ACM. ISBN 978-1-4503-3395-5/15/09\$15.00

DOI: <http://dx.doi.org/10.1145/2700171.2791059>.

International academic conferences have a diverse community, with different cultural backgrounds and languages. Thus, it is interesting to analyze how language affects the generation of content and interaction among attendees. Such study would allow to observe how integrated a research community is and will shed light for future research.

This can be of special interest to conference organizers not only to evaluate communication but also to have an overview of their audiences. Despite the research published in the past [6, 7, 13, 12] on academic conferences, little has been done on language communities and the communication established among them. To bridge this gap, we explore the language of 7M tweets posted by 18K users during 26 Computer Science conferences over five years (one week before and after for each conference).

We group users by the language(s) they use to tweet in order to explore how different language communities interact. Although English is expected to be the *lingua franca* of many international events, we wonder to what extent people use other languages on Twitter during academic conferences.

Research Questions. Overall, our study was driven by the following research questions:

- **RQ1. Conference attendees' languages:** To what extent do people tweet in other languages beyond English in conferences?
- **RQ2. Interactions between lingua groups:** How do lingua groups interact with each other?
- **RQ3. Effect of language:** Is there an effect of language or lingua group over online user interaction?

Main results. We find that most people tweet only in English (61%) in conferences but most of the tweets are posted by multilingual users and their participation varies significantly across conferences. Additionally, we observe that *English monolinguals* receive most of the attention and interact more within their lingua group while the opposite is observed for most of the members from other language communities.

Finally, we show that people who do not interact with other attendees are mostly monolinguals, while people who interact with others tweet in different languages (bilingual or trilingual).

2. DATASET

We selected a representative set of conferences in Computer and Information Science from the CORE Conference Ranking list¹; 26 conferences active in Twitter every year between 2009 and 2013. Furthermore, we manually checked that the selected conferences did not overlap with other events. To retrieve the tweets from these events in previous years, we used the Topsy API and crawled

¹<http://www.core.edu.au/index.php/conference-rankings>

Table 1: Percentage of monolinguals, bilinguals and multilinguals tweeting in each conference between 2009-2013 (col 2-4). Diversity percentage for interactions (reciprocated or not) (col 5-8) for tweets (TW) and retweets (RT).

Conference	Lingua groups			Diversity percentage			
				General		Reciprocated	
	1-ling	2-ling	≥ 3-ling	TW	RT	TW	RT
AAAI	81%	8%	11%	34%	29%	16%	20%
ACMMM	52%	38%	11%	53%	53%	48%	41%
CHI	76%	17%	7%	49%	48%	40%	30%
CIKM	66%	24%	10%	54%	54%	44%	40%
ECIR	58%	27%	15%	55%	57%	43%	31%
ECIS	57%	31%	12%	46%	44%	24%	0%
HT	64%	26%	10%	52%	53%	37%	29%
ICIS	67%	26%	7%	44%	41%	19%	16%
ICML	75%	17%	8%	52%	55%	20%	21%
ICMT	51%	30%	19%	70%	62%	31%	20%
ICSE	58%	32%	10%	47%	46%	40%	47%
ISMAR	64%	28%	8%	39%	37%	19%	21%
IUI	62%	21%	17%	59%	58%	45%	44%
KDD	73%	18%	9%	53%	50%	38%	37%
MobileHCI	66%	23%	11%	50%	47%	48%	39%
NIPS	74%	19%	7%	46%	48%	25%	20%
SIGGRAPH	77%	16%	7%	38%	32%	24%	19%
SIGIR	68%	21%	11%	56%	58%	36%	39%
SIGMOD	72%	23%	5%	58%	53%	19%	12%
SLE	59%	32%	9%	58%	58%	40%	40%
UBICOMP	71%	21%	8%	59%	57%	55%	44%
UIST	71%	24%	5%	60%	58%	35%	32%
VLDB	67%	26%	7%	56%	53%	29%	21%
WSDM	65%	22%	13%	61%	60%	48%	39%
WWW	52%	32%	16%	52%	51%	43%	40%
XP	58%	35%	7%	53%	52%	51%	54%

tweets containing the corresponding official hashtag (e.g., #chi12, #www2009) within a two-week time window around the dates each conference took place (from seven days before and until seven days after the conference ended). We found that these tweets were posted by 22,021 participants in total. We acknowledge that these participants also interact with others without the conference hashtag and because of this we also crawled their timeline tweets during the same period. In total, we obtained 6,993,693 tweets.

Language Identification. To identify the language of the tweets, we removed all URLs, mentions and hashtags. Then we set a minimum threshold of 4 remaining words in the tweets to identify their language. The language detection task was performed with a professional language tool provided by Yahoo! Labs that is able to identify over 40+ languages as in [9]. Following this process we were left with 6,184,775 tweets (88% from initial sample) with an identified language (see Table 5 in the Appendix). The tweets without a language are generally those containing symbols or links only. Finally, we proceeded to model each user by the three most frequent languages they used to tweet (setting a minimum threshold of 5 tweets per language). Consequently, we found 266 lingua groups with 18,347 users using at least three different languages in their tweets².

3. RESULTS

RQ1. To what extent do people tweet in other languages beyond English across conferences?

As expected, we found that the majority of tweets are written in English (76%). Nevertheless, due to the multicultural nature of conferences, there is a non-negligible 24% of tweets in languages different than English (en), such as French (fr), Spanish (es), German

²The venue of the conference (city, country) might have had an impact on the languages used, but we leave that detailed analysis for future work.

Table 2: Statistics of top lingua groups (more than 90 users). We show the percentage of users belonging to each *lingua* (Users), the percentage of tweets (Tweets), the engagement (tweets/user), and the interquartile range of tweets per users (IQR).

Lingua	Users	Tweets	(tweets/user)	IQR
en	61.31%	31.56%	167.18	142.00
en-fr	6.46%	3.85%	193.88	164.75
en-es	3.79%	2.52%	216.14	191.00
de-en	2.18%	1.75%	260.68	249.25
en-nl	2.15%	1.58%	238.65	222.5
fr	2.00%	0.27%	43.61	43.00
en-ja	1.92%	1.16%	196.3	166.00
en-es-pt	1.62%	4.05%	809.84	611.50
en-pt	1.44%	0.37%	83.52	63.75
en-it	1.36%	1.68%	402.87	193.50
nl	1.36%	0.16%	37.46	31.00
ja	1.09%	0.16%	47.94	42.75
en-es-fr	0.93%	9.23%	3,224.06	1,771.00
ca-en-es	0.79%	2.18%	891.29	799.50
en-ko	0.57%	0.53%	301.94	300.00
es	0.52%	0.06%	35.24	40.00
Others	10.52%	38.88%	1,200.51	864.00

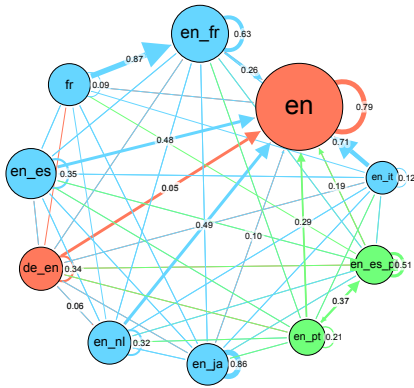
Table 3: *Most popular linguas*: lingua groups ordered by the attention they receive across all conferences. The *out-link* column represents the percentage of interactions going to other lingua groups.

General					
Mentions (148,184)			Retweets (91,523)		
Ling.	Att.	out-links	Ling.	Att.	out-links
en	67%	37%	en	66%	37%
en-fr	7%	56%	en-fr	7%	54%
de-en	3%	74%	de-en	3%	78%
en-es	3%	79%	en-es	3%	80%
en-ja	2%	35%	en-ja	2%	42%
Reciprocated					
Mentions (25,956)			Retweets (6,496)		
Ling.	Att.	out-links	Ling.	Att.	out-links
en	57%	48%	en	51%	52%
en-fr	8%	52%	en-fr	8%	44%
de-en	4%	72%	en-es	5%	61%
en-es	4%	71%	de-en	4%	74%
en-nl	3%	71%	en-nl	3%	70%

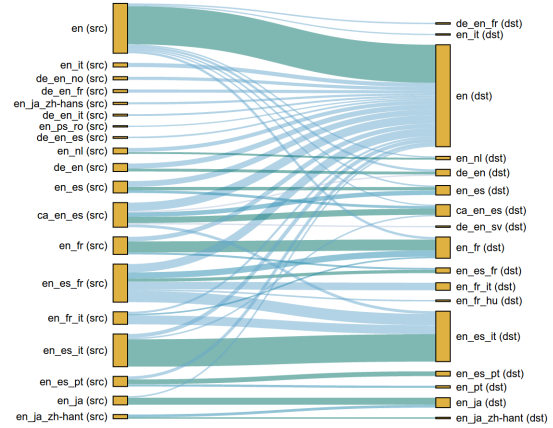
(de) and Japanese (jp). Furthermore, we found in our dataset that many people post tweets in more than a single language.

We quantify this observation in Table 1 that shows the percentage of users who tweet in a single language (1-lingua), in two languages (2-lingua) or three or more (≥ 3 -lingua) in each conference. We observe that the percentage of people who tweet in two or more languages goes from close to 20% (AAAI, SIGGRAPH) up to around 50% (ACMMM, ICMT, WWW) showing important differences among conferences in the distribution of users who tweet in one or more languages. Based on these results, rather than analyzing languages as isolated groups, we studied the lingua groups as communities of people who speak either one or more languages. Table 2 describes the top language communities by number of users. The table shows that the majority of users are classified as *English monolinguals* (61%) but interestingly only produce (29%) of all tweets with a moderate engagement (only 179.5 tweets per user). In contrast, we see that users of multilingual groups are the most engaged (3609.9 tweets/user for en-es-fr, 1016.7 for ca-en-es, and 944.93 for en-es-pt).

These results lead us to further analyze specific lingua groups to unveil the interaction between language communities and their online behaviour.



(a) Mentions between lingua groups. An edge from lingua x pointing to lingua y shows proportions of mentions that people in lingua x directed to people in lingua y . For readability, we only show probabilities ≥ 0.05 .



(b) Retweet interactions between top 50 most active lingua groups.

Figure 1: (a) Nodes representing the top 10 lingua groups based on mentions. (b) Interactions between lingua groups based on source language (src) retweeting posts in a target language (dst).

RQ2. How do lingua groups interact with each other?

To answer this question, we first define two types of interactions: (1) general interactions and (2) reciprocated interactions. We refer to *general interactions* to all (re)tweets containing mentions in the original text of the tweet, while *reciprocated interactions* correspond to the (re)tweets that were reciprocated by the users mentioned in the text.

Secondly, we measure diversity using the Gini-Simpson index, as in [3, 5] also called *diversity index*. This index ranges from 0 to 1 and it measures the probability that two lingua groups taken at random from a set of interactions are different. Participants of a conference with diversity index close to 0 will have the tendency to interact with people of the same lingua group. Conversely, values close to 1 show a uniform distribution of interactions with other lingua groups. We define diversity D of a lingua group as:

$$D(c, i) = 1 - \sum_{j \in S} \left(\frac{I_{i,j}^c}{N_i} \right)^2 \quad (1)$$

with $N_i = \sum_{k \in S} I_{i,k}^c$ and where $I_{i,j}^c$ is the total number of interactions between people of lingua i and j . N_i is the total number of interactions of people of lingua i in conference c . To know the diversity of a conference, we average $D(c, i)$ over all the linguas in conference c .

We see in Table 1 the diversity percentage for each conference. We find some interesting patterns showing that a lower percentage of monolinguals is linked to higher diversity. For example, ICMT is the most diverse conference for the general type of interactions and the percentage of monolinguals is the lowest of all (51%). Conversely, AAI shows high percentage of monolinguals (82%) and the lowest diversity for the general interactions. On the other hand, reciprocal interactions do not show to be related to the percentage of monolinguals. For example, UBICOMP presents a high percentage of monolinguals and the highest diversity for the reciprocal interactions.

Furthermore, we look at the attention *received* by members of each lingua by calculating the number of mentions and retweets received from different users. Table 3 shows the top 5 most popular lingua groups: English monolinguals are the most mentioned and retweeted in the general and reciprocated interactions. Albeit the fact that English monolinguals do not produce most of the tweets, they still receive most of the attention. This is mostly explained by the column *out-links*, which shows the percentage of mentions and retweets about *different* lingua group. For example, we see that only

37% of the mentions and retweets generated by English monolinguals refer to other groups. Interestingly, Japanese bilinguals also prefer to interact mostly within their group. Conversely, groups like *en-fr*, *de-en*, *en-es* interact more with users of *different* lingua groups.

The unequal activity between lingua groups is also seen in Figure 1, which considers only the top 10 lingua groups and shows (a) the mentions network (general type) and (b) the retweet network (general type) across lingua groups. Figure 1a shows that 79% of all mentions from the *en* group also belong to the same group. Moreover, 35% of mentions from the *en-es* lingua group refer to users from the same group, and 48% to the *en* group.

In Figure 1b, the Sankey plot represents the network of retweets. Again, here we see that for most of the cases the English group retweets members from the same group. At the same time, the English group receives most of the attention from other language communities. Interestingly, in similar proportion, lingua groups *en-es-it*, *en-fr*, *en-es-pt* and *en-ja* show a similar pattern, preferably retweeting users on their same lingua groups.

RQ3. Is there any effect of language or lingua group over online user interaction?

We addressed this question by studying how the number of languages a Twitter user speaks affects her online behavior. As already explained, if a user has posted tweets in only one language we consider her in the 1-lingua group (monolingual), while another user tweeting in two languages will be in the 2-lingua group, and so on. We found two results that show at general and at individual level the effect of the amount of languages on user interaction. At the general level, we found that among the users who posted tweets but who had not interacted with other people (by mentioning them), the percentage on monolinguals is considerably larger (80.6%) than multilinguals. A different picture is seen among users who interacted at least once during the conference (by mentioning someone in a tweet), since only 62.9% of those users are monolinguals and the rest are multilinguals. We conducted a chi-square test of proportions comparing the distribution of monolinguals, bilinguals and trilinguals between people who interacted (using mention or retweet) and people who did not. We found a statistically significant difference with $\chi^2 = 416.6$, $df = 2$, $p < .001$. This relation can be better observed in Figure 2, where the group who interacted (right-side plot) had a more balanced distribution and hence a higher entropy (a measure of diversity [11]) of $H(s) = 0.89$ compared to a smaller diversity on lingua groups among people who did not interact with

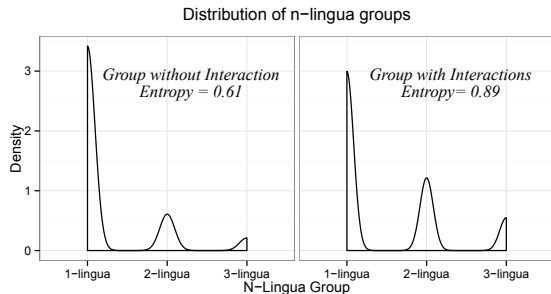


Figure 2: Distribution of users per n-lingua groups: users who interacted (right) and did not interacted (left) with others.

an entropy $H(s) = 0.61$. Moreover, at the individual level we found that the more the languages a user speaks, the larger the likelihood to interact with others. Table 4 shows the results of a logistic regression where the dependent variable measures whether the user *interacted* with other people or not. The factors in the regression are the *year* of the conference and the number of languages the user has used to tweet (*n_languages*). We observe that the number of languages has a significant β coefficient of 0.666 ($p < .001$), which can be interpreted by saying that, keeping all the other factors fixed, for each additional language the user speaks the odds ratio of interacting in the network increases by 95% (since $e^{0.666} = 1.95$).

4. RELATED WORK

The role of Twitter in academic conferences was also studied by Letierce *et al.* [6, 7]. They showed that Twitter is frequently used to spread information using the official conference hashtags. Wen *et al.* [13, 12] found that newcomer students receive little attention from senior members of the research community and identified factors that contribute to the continuing participation of users to the online Twitter conference activity. We have continued this line of research by exploring the influence of language during conferences. The role of language in Twitter has also been studied. Hong *et al.* [4] studied differences in usage patterns between language communities in Twitter, while Kim *et al.* [5] performed a sociolinguistic study on the role of mono- and bilinguals in Twitter across multilingual societies such as Qatar, Quebec and Switzerland. Inspired by them, we adopted similar methods to build language communities but we targeted different lingua groups interacting at conferences.

A broader but certainly related topic of study is the impact of *culture* in online communication. Garcia *et al.* [2] found that language and cultural dimensions are discriminative features influencing international active conversation and attention in Twitter. We find that focusing on language(s) we capture the multicultural nature of most researchers that attend international conferences.

5. CONCLUSIONS & FUTURE WORK

In this paper we show that most of the English tweets posted in Computer and Information Science conferences come from lingua communities different than English monolinguals. We also observe that English monolinguals still prefer to interact more with themselves. The same happens for other important communities such as English-Japanese bilinguals, while this behavior is not commonly replicated in other communities, who tend to interact more equally with members of other linguas.

We also find that there is more language diversity among people who interact with others on Twitter during conferences, compared to people who do not. This result suggests an important implication, which is that although English is the standard for scientific communication, the diversity in language use is a catalyst for interactions in a community.

Table 4: Results of L.R. where the D.V. is whether user interacted on Twitter (mentions) and the I.V.s are conference *year* and number of languages spoken.

Variable	β coeff.	S.E.
year(=2009)	2.049***	(0.390)
year(=2010)	2.458***	(0.385)
year(=2011)	2.453***	(0.385)
year(=2012)	2.294***	(0.383)
year(=2013)	2.423***	(0.383)
n_languages	0.666***	(0.035)
Constant	-1.371***	(0.385)
Observations	26,281	

Note: *p<0.1; **p<0.05; ***p<0.01

There are still several questions to address in future work. For example, does the attention received by the English monolingual group depend on popular users and/or venues? Furthermore, we will validate the result of RQ3 using a “test set” of tweets in the future.

Acknowledgments: This work was carried out during the tenure of an ERCIM “Alain Bensoussan” fellowship program by C.T.

6. REFERENCES

- [1] M. Ebner. Introducing live microblogging: How single presentations can be enhanced by the mass. *Journal of research in innovative teaching*, 2009.
- [2] R. García-Gavilanes, Y. Mejova, and D. Quercia. Twitter Ain’T Without Frontiers: Economic, Social, and Cultural Boundaries in International Communication. In *Proc. CSCW’14*, 2014.
- [3] R. Guimera and L. A. N. Amaral. Functional cartography of complex metabolic networks. *Nature*, (7028), feb 2005.
- [4] L. Hong, G. Convertino, and E. H. Chi. Language matters in twitter: A large scale study. In *Proc. ICWSM’11*, 2011.
- [5] S. Kim, I. Weber, L. Wei, and A. Oh. Sociolinguistic analysis of twitter in multilingual societies. In *Proc. HT’14*, 2014.
- [6] J. Letierce, A. Passant, J. Breslin, and S. Decker. Understanding how twitter is used to widely spread scientific messages. In *Proc. WebSci’10*, 2010.
- [7] J. Letierce, A. Passant, J. G. Breslin, and S. Decker. Using twitter during an academic conference: The #iswc2009 use-case. In *Proc. ICWSM’10*, 2010.
- [8] E. Mitchell and S. B. Watstein. The places where students and scholars work, collaborate, share and plan: endless possibilities for us! *Reference services review*, 35(4), 2007.
- [9] B. Poblete, R. García-Gavilanes, M. Mendoza, and A. Jaimes. Do All Birds Tweet the Same? Characterizing Twitter Around the World. In *Proc. CIKM’11*, 2011.
- [10] C. Ross, M. Terras, C. Warwick, and A. Welsh. Enabled backchannel: conference twitter use by digital humanists. *Journal of Documentation*, 67(2), 2011.
- [11] C. E. Shannon. A mathematical theory of communication. *ACM SIGMOBILE Mobile Computing and Communications Review*, 5(1):3–55, 2001.
- [12] X. Wen, Y. Lin, C. Trattner, and D. Parra. Twitter in academic conferences: Usage, networking and participation over time. In *Proc. HT’14*, 2014.
- [13] X. Wen, D. Parra, and C. Trattner. How groups of people interact with each other on twitter during academic conferences. In *Proc. CSCW’14*, 2014.

APPENDIX

The following tables show detailed data used in our analyses.

