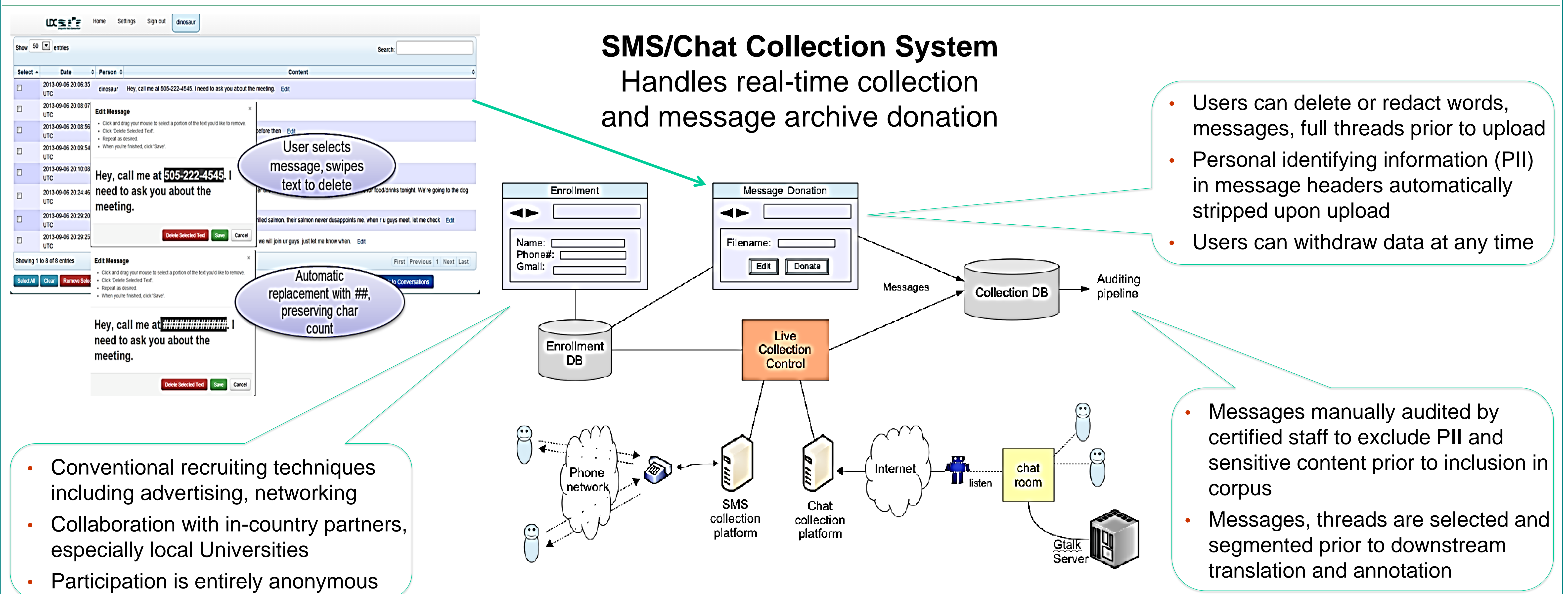


Collecting Natural SMS & Chat Conversations in Multiple Languages: BOLT Phase 2 Corpus

Z. Song, S. Strassel, H. Lee, K. Walker, J. Wright, J. Garland, D. Fore, B. Gainor, P. Cabe, T. Thomas, B. Callahan, A. Sawyer

- **The Broad Operational Language Translation (BOLT) Program** is developing technology that enables English speakers to retrieve and understand information from informal foreign language sources including chat, text messaging (SMS) and spoken conversation
- **BOLT Phase 2 SMS/Chat Collection Requirements**
 - **Multilingual:** English, Chinese, Egyptian Arabic
 - **Large volume:** at least 2 million words/language
 - **Conversational:** naturally-occurring, two-way, informal SMS and chat messages
- **Prior SMS/Chat Corpora** are primarily one-sided conversations or public chat room transcripts with limited content



Collection Result

Language	Active Users	Collect Words	Audit Pass Rate	Final Words	Final Msgs	Final Convs	Avg Messages/Conv	Avg Words/Message	Avg Words/Conv	Avg Words/User	Collection Method		Genre	
											Live	Donation	SMS	Chat
Egyptian	26	690K	69%	475K	119,00	2140	56	4	222	48,051	4%	96%	39%	61%
Chinese	77	3.7M	54%	2M	306,99	7844	39	7	255	21,710	1%	99%	4%	96%
English	152	3.3M	79%	2.6M	212,000	9155	23	12	284	27,600	6%	94%	94%	6%
Total/Average	255	7.69M	67%	5.075M	212,000	19,139	39.33	7.67	253.67	32453.67				

Linguistic Features of Collected Data

English

Argh I suck at txt srry

I know I still need to give you your hoodie back

Yes u def suck at text

Lol I'm used to it now tho

Hahah good

*good you don't take it personally

When do u think Ur moving to conshy?

ldk yet :(

Chinese

你最近回来纽约吗? 小周的茶叶还在这里呢

会来

最近没机会 月底要出长差

去哪里? 什么时候回?

去冰岛 五月24日去 六月一日会

美国有没房屋租赁的网站? 因为这次就我一个人, 单位嫌以前的太大, 需要换个单室套

不是特别清楚 最好问本地人

Arabic

Enty ya benty msh btrodii 3la elbta3 da abdan

el net kan mtnayel

eh akhbar el sa23a 3ndkim

Tab matsha3'ali 5edmt el iphone men Vodafone

eh da nezmha eh de??!!

Bs 7atsh3'lk kol 7aga

بجد؟ زي ايه يعني؟

Elnet w el apps

General Features

- Natural and conversational
- Internet slang and novel abbreviations
- Nonstandard use or complete absence of punctuation
- Frequent typographical errors, misspellings and missing spaces
- Widespread use of simple and complex initialisms
- Prevalent pro-drop

Arabic Features

- Code mixing, neologisms
- Multiple orthographies
 - 63% Arabizi
 - 13% Arabic
 - 24% Mixed
- Arabizi data semi-automatically transliterated into Arabic script for use in BOLT
- Joint work with Arabic Dialect Group at Columbia University