

学会報告 IJCNLP-AAACL 2023 (第13回自然言語処理国際共同学 会・第3回計算言語学会アジア 太平洋支部大会)

野元 裕樹 (nomoto@tufs.ac.jp)
Luncheon Linguistics, 2023/11/29



この報告で伝えたいこと

1. 学会の構成
2. 個々の発表の紹介
言語記録活動 (language documentation) や言語資源開発 (language resource development) に関連するもの
3. 言語学の学会との違い
自然言語処理・計算言語学の学会の国際学会に初めて参加してみて知ったこと

3

学会の概要

- The 13th International Joint Conference on Natural Language Processing (IJCNLP)
- The 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (AAACL)
 - Cf. North American Chapter (NAACL), European Chapter (EACL)
- 日程：2023年11月1-4日
- 場所：インドネシアバリ島
グランドハイアットホテル
- 実施形態：ハイブリッド（ポスターは会場のみ）
- 大会HP：<http://www.ijcnlp-aacl2023.org>



2

学会の構成

- 1日目：ワークショップ、チュートリアル
- 2-4日目：本大会
 - Main tracks (oral, poster, findings*)
 - Keynote speeches
 - Featured plenary talks
 - System demonstrations

Findings

- ここ最近ACLで始まったものらしい
- 口頭発表・ポスター発表での選択には至らなかった論文
- 予稿集には掲載される
- Findingsのセッションでは一発表2-3分で次々と論文の紹介をしていた

4

ワークショップ

- ・ 応募して採択されたもの
- ・ 採択されたワークショップはそれぞれ発表者の募集・選抜を行う

1. Student Research Workshop
2. 3rd Workshop on NLP for Medical Conversations
3. Workshop on Information Extraction from Scientific Publications
4. The Fourth Workshop on Evaluation & Comparison of NLP Systems
5. The 11th International Workshop on Natural Language Processing for Social Media
6. **The First Workshop on South East Asian Language Processing**
 - ・ 「東南アジア」とはいうものの、インドネシアとフィリピンのみ。
 - ・ 南インドの発表が2件...
7. The 6th Workshop on Financial Technology and Natural Language Processing
8. Second Workshop on Natural Language Interfaces
9. The ART of Safety: Workshop on Adversarial testing and Red-Teaming for generative AI

5

チュートリアル

- ・ 応募して採択されたもの
- ・ 日本言語学会の「ワークショップ」に相当

1. Language and Robotics: Toward Building Robots Coexisting with Human Society Using Language Interface
2. **Current Status of NLP in South East Asia with Insights from Multilingualism and Language Diversity**
3. Practical Tools from Domain Adaptation for Designing Inclusive, Equitable, and Robust Generative AI
4. Editing Large Language Models
5. Learning WHO Saying WHAT to WHOM in Multi-Party Conversations
6. Developing State-Of-The-Art Massively Multilingual Machine Translation Systems for Related Languages

6

本大会プログラム

詳細はハンドブックを参照
https://github.com/IJCNLP-AAACL23-Files/handbook/blob/main/IJCNLP_AAACL_Handbook_2023_v5.pdf

1日目	Dialog Systems and Generation	Question Answering
	Demo	Regional Language Processing
2日目	Demo	Findings Session
	Resources and Evaluation	Data Mining, Information Extraction and Retrieval
	Poster Session	
3日目	Multilingual and Multimodal Analysis	Machine Learning and Model Interpretability
	Semantics	NLP Applications

7

8

個々の発表の紹介



1. 私の発表（ポスター）

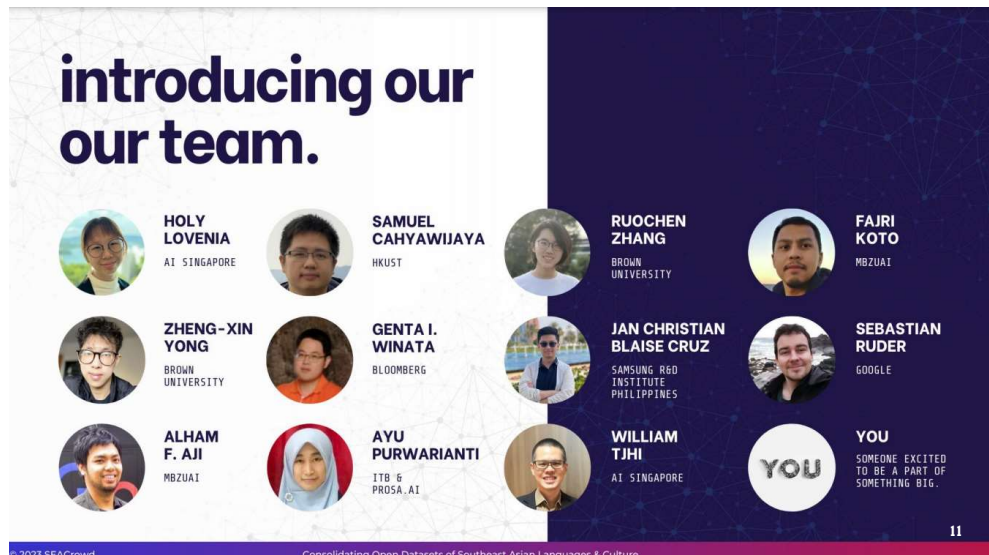
Issues Surrounding the Use of ChatGPT in Similar Languages: The Case of Malay and Indonesian

- ChatGPTにマレー語で聞くとインドネシア語で返ってくる問題
 - 実際の頻度
 - 問題が起きる理由
 - 問題をもたらす社会的・倫理的問題
 - 問題を解決する技術的・社会政治的方法
- ポスター実物は野元研究室（603）前に掲示中
<http://www.tufs.ac.jp/ts/personal/nomoto/IJCNLP-AAACL2023-poster.pdf>
- 予稿集論文
<http://www.afnlp.org/conferences/ijcnlp2023/proceedings/main-short/cdrom/pdf/2023.ijcnlp-short.9.pdf>



2. SEA CROWD（チュートリアルから）

- 東南アジアの言語・文化の公開されたAI開発用データセットを一か所に集め、かつ標準化する試み
- チュートリアルスライド pp.107-
https://aacl2023-sea-nlp.github.io/assets/tutorial_slides.pdf



SEACROWDへの貢献方法

- 既存の公開データセットのメタデータを送る
- #1のデータセットのためのデータローダーを作成する
- （個人のPCや大学研究室のサーバーにあるとされる）未公開のデータセットの情報を提供する
- 自らの未公開のデータセットを公開する（+ #3の情報提供）

3. NUSAWRITES

- 予稿集論文: <http://www.aflnp.org/conferences/ijcnlp2023/proceedings/main-long/cdrom/pdf/2023.ijcnlp-long.60.pdf>

NusaWrites: Constructing High-Quality Corpora for Underrepresented and Extremely Low-Resource Languages

Samuel Cahyawijaya^{1,4*}, Holy Lovenia^{1,4*}, Fajri Koto^{3,4*}, Dea Adhista^{2†}, Emmanuel Dave^{4†}, Sarah Oktavianti^{2†}, Salsabil Maulana Akbar^{4†}, Jhonson Lee^{4†}, Nuur Shadieq^{4†}, Tjeng Wawan Cenggoro^{8,4†}, Hanung Wahyuning Linuwih^{2†}, Bryan Willie^{1,4†}, Galih Pradipta Muridan^{2†}, Genta Indra Winata^{5,4†}, David Moeljadi^{7†}, Alham Fikri Aji^{3,4†}, Ayu Purwarianti^{6,4}, Pascale Fung¹

¹HKUST ²Prosa.ai ³MBZUAI ⁴IndoNLP ⁵Bloomberg
⁶Institut Teknologi Bandung ⁷Kanda University of International Studies ⁸Bina Nusantara University
 *Equal Contribution † Equal Contribution

3. NUSAWRITES

- 超低資源言語のコーパス作成は主にウェブスクレイピングや翻訳により行われてきた。
 - 利点: ①効率的②低コスト
 - 欠点: ①語彙的多様性がない②現地の文化と無関係な内容
- インドネシアの地方語を対象に母語話者による翻訳 (NusaTranslation) とパラグラフライティング (NusaParagraph) のラベル付きコーパスを構築し、Wikipediaデータと比較して、欠点が解消されることを示した。
- 新たなベンチマーク (評価指標) NUSAWritesの構築・利用
 - 言語理解タスク (感情認識、感情分析、談話のモードの分類、トピックモデル)
 - 言語生成タスク (機械翻訳)
- データ・コード公開ページ: <https://github.com/IndoNLP/nusa-writes>

Language	ISO code	Status
Amboinese Malay	abs	wider communication
Mandailing	btm	threatened
Batak Tobu	btc	threatened
Betawi	bew	threatened
Bima	bhp	vigorous
Buginese	bug	wider communication
Javanese	jav	educational
Madurese	mad	developing
Makassarese	mak	threatened
Minangkabau	min	developing
Palembang / Musi	mui	wider communication
Rejang	rej	vigorous
Sundanese	sun	developing

4. MASAKHANEWS

- 予稿集論文: <http://www.aflnp.org/conferences/ijcnlp2023/proceedings/main-long/cdrom/pdf/2023.ijcnlp-long.10.pdf>

MasakhaNEWS: News Topic Classification for African languages

David Hesolwa Adelan¹, Murek Mstak¹, Israel Abebe Azime², Amelha Oluwalara Akabi³, Amelha Lashbea Tsegay³, Christine Mwangi⁴, Oluwasegun Ogunmola⁵, Isaacson E. P. Dossou^{6,7,8}, Akintunde Oluwalope⁹, Dorcas Ntindori, Chris Chimwene Emraz¹⁰, Sara Salah al-sazzani¹¹, Blessing K. Sikanda, Davis David¹², Lovethin Nidebe, Jonathan Mubili¹³, Iande G. Ajayi¹⁴, Tefana Mwan Ngazi¹⁵, Brian Okilombe, Abimbola Toluwalade Osofisan, Oluwaseun C. Oluwalara, Mahdiu Mohamed¹⁶, Shamsudeen Hassan Muhammad¹⁷, Teshome Mulugeta Akhalu¹⁸, Saleel S. Abdulkhalil¹⁹, Mwangi Genneth Yikraz²⁰, Tahmiddeen Goralabo, Hiba Abdullhameed²¹, Mahdi Taye Bane, Oluwalaseun O. Adegbenro²², Isamulhwa Shide²³, Tolope Amu Adedani, Habiba Abdalgany Kallani, Abdu-Rakeman Omosayo²⁴, Adetola Adekun, Adedabi Abereh, Amulawogbo Ayemba Oluwalaseun Summi²⁵, Chikwendia Sini²⁶, Waseem Khatib²⁷, Oyekareli Raphael Ogbe, Chinedu E. Mbonu²⁸, Chiamaka I. Chukwura²⁹, Samuel Eantij³⁰, Jessica Ojo, Oyinokunoluwa E. Awusan, Tadese Kebede Gugge³¹, Sakayo Indoum Sarf³², Pascha Nyirakye, Ewendence Ndikumana³³, Orono Vissouf, Mawliyah Oluwalara³⁴, Karama El Tokim, Usen Kimani³⁵, Thina Diko, Snyanda Nsohama, Sinodot G. Ngunjiri³⁶, Abdulmajid Tind Johari, Shafiq Abul Mohamud³⁷, Foad Miro Hassan³⁸, Mages Ahmad Mohamed³⁹, Edward Ngalire⁴⁰, John Drognyana, Ivan Norkovska, and Petrus Steynberg⁴¹

¹Mathematics NAE Africa, ²University College London, United Kingdom, ³Stanford University, Germany, ⁴Indiana Polytechnic National, Mexico, ⁵Yale University, China, ⁶University of Waterloo, Canada, ⁷Leipzig AI, ⁸McGill University, Canada, ⁹Mila Quebec AI Institute, Canada, ¹⁰Laifika, ¹¹Toronto University of Munich, Germany, ¹²Leipzig University of Technology, Germany, ¹³University of Lagos, Nigeria, ¹⁴University of Lagos, Nigeria, ¹⁵University of Lagos, Nigeria, ¹⁶University of Lagos, Nigeria, ¹⁷University of Lagos, Nigeria, ¹⁸University of Lagos, Nigeria, ¹⁹University of Lagos, Nigeria, ²⁰University of Lagos, Nigeria, ²¹University of Lagos, Nigeria, ²²University of Lagos, Nigeria, ²³University of Lagos, Nigeria, ²⁴University of Lagos, Nigeria, ²⁵University of Lagos, Nigeria, ²⁶University of Lagos, Nigeria, ²⁷University of Lagos, Nigeria, ²⁸University of Lagos, Nigeria, ²⁹University of Lagos, Nigeria, ³⁰University of Lagos, Nigeria, ³¹University of Lagos, Nigeria, ³²University of Lagos, Nigeria, ³³University of Lagos, Nigeria, ³⁴University of Lagos, Nigeria, ³⁵University of Lagos, Nigeria, ³⁶University of Lagos, Nigeria, ³⁷University of Lagos, Nigeria, ³⁸University of Lagos, Nigeria, ³⁹University of Lagos, Nigeria, ⁴⁰University of Lagos, Nigeria, ⁴¹University of Lagos, Nigeria



4. MASAKHANEWS

- アフリカの16言語のニューストピック分類用データセット
- データ・コード公開ページ: <https://github.com/masakhane-io/masakhane-news>

Language	Family/branch	Region	# speakers	News Source	# articles
Amharic (amh)	Afro-Asiatic / Ethio-Semitic	East Africa	57M	BBC	8,204
English (eng)	Indo-European / Germanic	Across Africa	1268M	BBC	5,073
French (fra)	Indo-European / Romance	Across Africa	277M	BBC	5,683
Hausa (hau)	Afro-Asiatic / Chadic	West Africa	77M	BBC	6,965
Igbo (ibo)	Niger-Congo / Volta-Niger	West Africa	31M	BBC	4,628
Lingala (lin)	Niger-Congo / Bantu	Central Africa	40M	VOA	2,022
Luganda (lug)	Niger-Congo / Bantu	Central Africa	11M	Gambuzee	2,621
Najja (pcm)	English Creole	West Africa	121M	BBC	7,783
Oromo (orm)	Afro-Asiatic / Cushitic	East Africa	37M	BBC	7,782
Rundi (run)	Niger-Congo / Bantu	East Africa	11M	BBC	2,995
chiShona (sna)	Niger-Congo / Bantu	Southern Africa	11M	VOA & Kwayedza	11,146
Somali (som)	Afro-Asiatic / Cushitic	East Africa	22M	BBC	2,915
Kiswahili (swa)	Niger-Congo / Bantu	East & Central Africa	71M-106M	BBC	6,431
Tigrinya (tig)	Afro-Asiatic / Ethio-Semitic	East Africa	9M	BBC	4,372
isiXhosa (xho)	Niger-Congo / Bantu	Southern Africa	19M	Isolozwe	24,658
Yorubi (yor)	Niger-Congo / Volta-Niger	West Africa	46M	BBC	6,974

言語学の学会との違い



言語学の学会との違い①

- 参加費が高い！
- 若い人たちは指導教員や勤務先が支払うのが一般的(？)

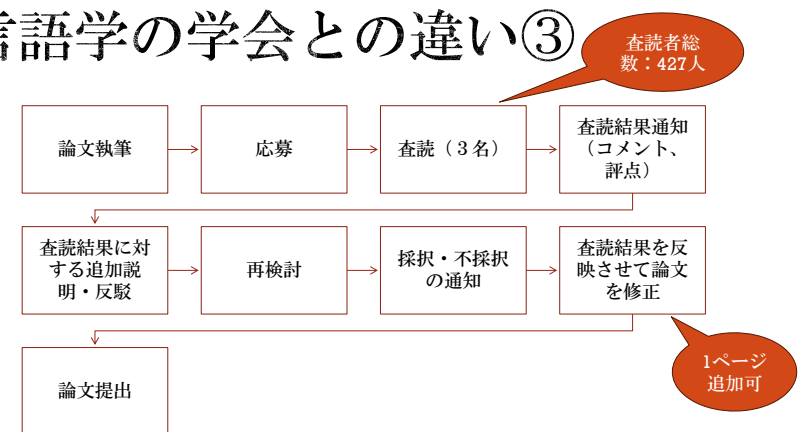
Timeline	Early Bird (by 30 Sep 2023)	Standard (by 21 Oct 2023)	Onsite (After 21 Oct 2023)	Remark
Regular (Normal)	750	850	950	Normal rate for non-student
Student (Normal)	450	550	650	Normal rate for fulltime student
Regular (Special)	700	800	900	Not to cover a paper. Special rate is for non-student from developing countries/regions.
Student (Special)	400	500	600	Not to cover a paper. Special rate for fulltime student from developing countries/regions.
Regular (Online)	300	400	500	Online participation, not to cover a paper
Student (Online)	200	225	250	Online participation, not to cover a paper

言語学の学会との違い②

- 応募先：直接 or ARR (ACL Rolling Review) という学会共通応募システム <https://aclrollingreview.org>
- 種目：Long Paper (8ページ+参考文献・付録) or Short Paper (4ページ+参考文献・付録)
- 発表形態（口頭orポスター）は内容に応じて主催者が決定

	Long			Short			計
	直接	ARR	Finding	直接	ARR	Finding	
応募数	229	19	-	109	6	-	363
採択数	72	23		22	14		131
採択率	29.0%	(38.3%)		19.1%	(31.3%)		36.1%

言語学の学会との違い③



言語学の学会との違い④

- Limitations、Ethics Statementの節が論文中で必須（ページ制限外）
Ethics Statementの例（NusaWritesの論文より抜粋、強調は発表者による）
 - Within this work, the annotators are properly rewarded above the national average minimum wage in Indonesia.
 - We have obtained informed consent from all annotators and adhered to data protection and privacy regulations for releasing the corpus and benchmark.
 - Throughout our research process, we have made conscious efforts to engage with the language communities, involve local experts, and respect their linguistic and cultural nuances.
- Responsible NLP Checklistを論文と一緒に提出
<https://aclrollingreview.org/responsibleNLPresearch/>

21

言語学の学会との違い⑤

- 参加者が若い
Are you a Ph.D. student?
- 中国、インド、中東からの参加者が多い（逆に、日本からの参加者が少ない）
- 企業からの参加者もいる
Do you work in industry or academia?
- スポンサー企業がいる
- 賞がたくさんある
- 口頭発表でNo showが結構ある

22