

Developing the EMI Corpus of Student Academic Reading & Writing: Addressing challenges in a large-scale, transnational collaborative project

Dana Gablasova - Raffaella Bottini - Vaclav Brezina - Luke Harding - Sally Ren Lancaster University





# Overview of this session

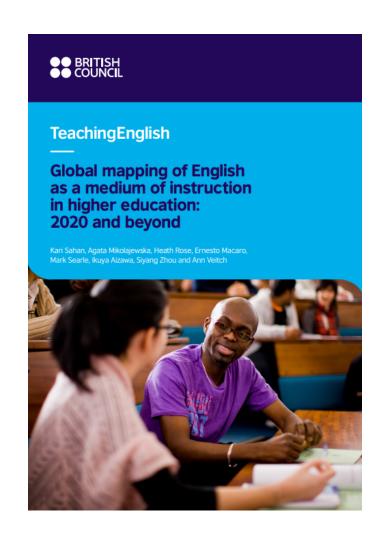
1. EMI Corpus project: introduction

method suzuki: Breach Bo with arboxilic acid

- 2. Overview of the EMI corpus
- 3. Challenges in the EMI Corpus multi-site collaborative project

VStept method. \_ carbo cation reach

# English medium instruction (EMI)



- EMI teaching/learning disciplinary subjects through the medium of English typically in countries where English is not the community language (McKinley, 2024; Pecorari & Malmström, 2018).
- EMI currently a global pedagogical trend; on the increase
- Use and knowledge of English crucial for understanding subject knowledge and for learning

# EMI: Challenges related to language use

- We know a lot about EMI reported via surveys, interviews, classroom observations, document analysis
- We know that students report difficulties related to speaking, writing and reading English – with potentially negative consequences for their academic success
- However, we do not have much data about how they actually use English and what demands are placed on them (e.g. in their reading) → calls for corpus research in EMI (Jablonkai, 2021)

## EMI Corpus project





Project: "Linguistic demands of EMI in Higher Education: A corpus-based analysis of student writing and reading in EMI university settings."

Funded by the British Council as part of the **Future of English research scheme** for 2022-25



Prince of Songkla University



Thammasat University



**University of Turin** 



University of Milan



Xi'an Jiaotong University



Xi'an Jiaotong-Liverpool University



Vienna University of Economics and Business

Corpora of EMI language use



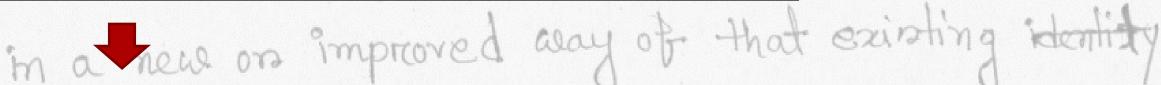
# Corpus evidence and EMI

oceating a me product,

Description of linguistic patterns and regularities

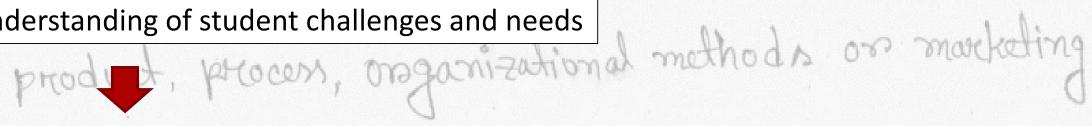


Understanding what language students produce and encounter



Understanding of student challenges and needs

Amourn to the question no: 3



Inform language teaching and testing practice/materials; Inform EMI policy (e.g., admission requirements, EAP provision, ESP provision)

magnizational methods or marketing methods

MA International Relations - MSc Computer Science — MBA - MSc Project Management - MA Intercultural Communication - MEng Mechanical Engineering - MA History — MBA - MA History - MSc Business Analytics - MSc Advanced Marketing Management - MEng Hons Chemical Engineering - MSc Finance - MSc Developmental Psychology - MSc Digital Business, Innovation and Management - MSc Human Resources Management - MA International and Military History - MSc Computing - MSc Mechanical Engineering - MA Creative Writing with English Literature - MA in History - PGCert Regional and Local History - MA Philosophy and Religion - BEng Chemical Engineering - MSc in Engineering - MA Language and Linguistics- BEng

# Corpus of EMI reading and writing: Overview

MA in Leadership and Management - MA Management Science - MSc Developmental Disorders - MA Language and Linguistics - MSc Digital Business Innovation and Management - BEng Hons Mechatronic Engineering - MSc Advanced Mechanical Engineering - MA Applied Linguistics - MA Applied Linguistics and TESOL - MSci Biomedicine - MSc Entrepreneurship and Innovation - MA Politics - MA Digital Humanities - MSc Criminology and Social Research Methods - MA Digital Humanities - MEng Mechanical Engineering - MA Linguistics - MSc Biomedicine - MSc Advanced Marketing Management - MChem Chemistry Hons - MEng Mechanical Engineering - MA Media and Cultural Studies - MSc Business Analytics - MA English Literary Studies - MA Creative Writing (Distance Learning) - MA Criminology and Criminal Justice - MA Discourse Studies - MSc Advanced Mechanical Engineering - MA Social Justice and Education - MA Philosophy - MSc Health Research - MA Intercultural Communication - MSc Conservation and Biodiversity - MSc Volcanology and the Environment

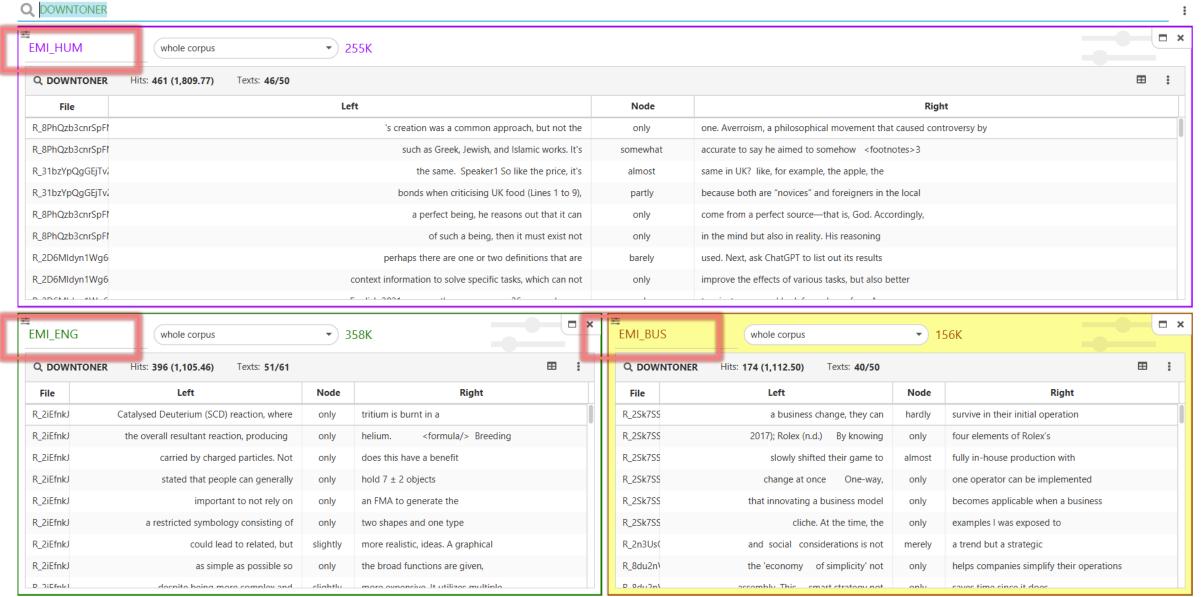
# Corpus size: 4.5M words and 2,000+ student papers

Level	Undergraduate		Postgraduat	е	
Disciplinary area	Business & Management		Humanities & Soc. Science	Science & Technology	
Core subjects	Business studies, Management, Finance, Accounting,		History, Literature, Sociology, Linguistics	Chemistry, Physics, Engineering	
Corpus size (Italy, Thailand, China, Austria)	234	283	525	345	
Corpus size (UK-based data)	20	140	257	200	
Corpus size (data total)	254	423		545	

# Frequency information

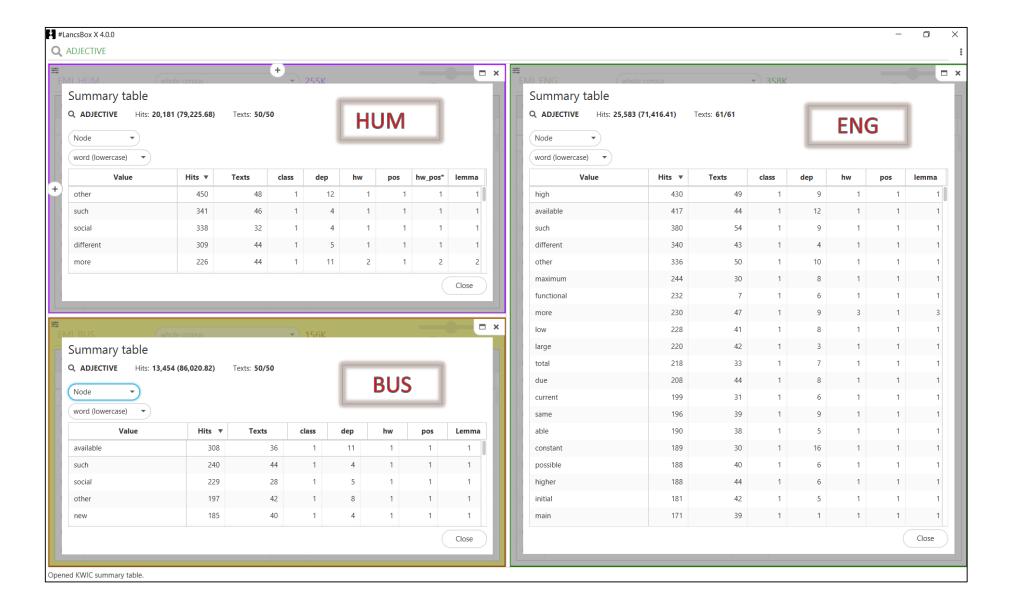
#### #LancsBox X 4.0.0

#### What downtoners are used in different disciplinary areas?



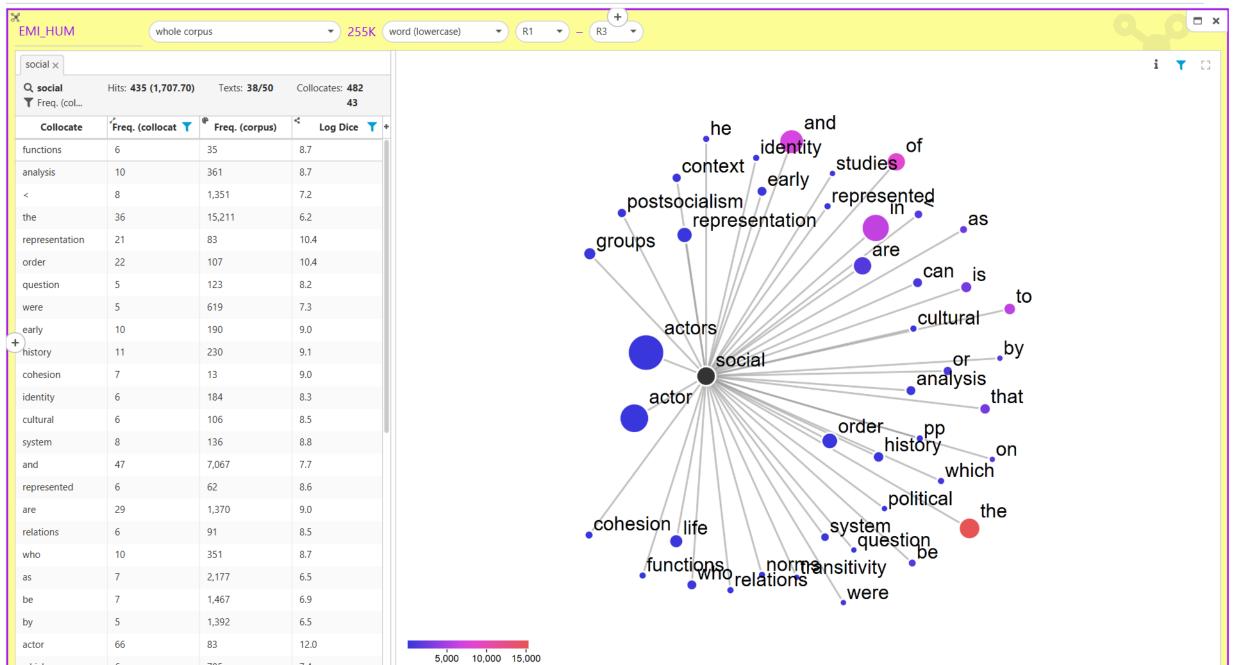
C----L--- IMMIC 4-- "DOMNITONIED"

#### What adjectives are common in different disciplines?



# Collocation analysis: Collocations

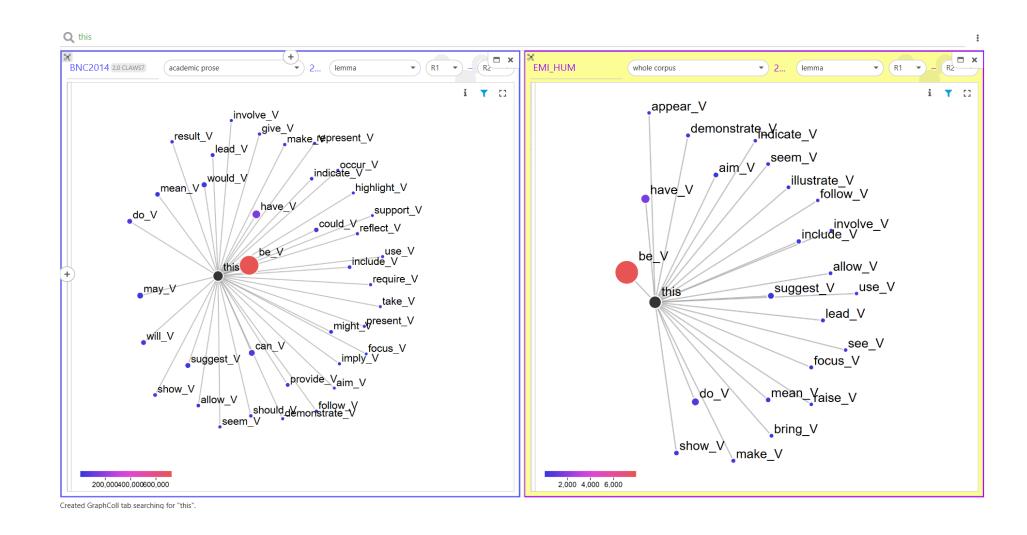
What do you think would be the strongest collocate for the word "social" in the writing of students in the humanities?



E	NG		HUM		BUS	
1. re	port	129	essay	85	essay	33
2. pr	oject	120	study	49	approach	34
3. va	lue	75	paper	33	report	21
4. pr	ocess	57	test	41	case	15
5. st	udy	49	case	30	study	13
6. ca	ise	47	way	26	method	13
7. se	ection	44	period	23	model	15
8. m	ethod	42	section	21	variable	13
9. m	odel	28	dissertation	16	strategy	14
10. de	esign	26	approach	17	perspective	12
11. w	eek	23	research	19	project	18
12. pc	oint	21	time	15	section	10
13. sy	stem	21	point	13		
14. re	search	21	image	12		
15. in	vestigation	20	task	12		
16. da	atum	20	listening	11		
17. ed	quation	19	context	11		
18. re	sult	19	perspective	10		
19. to	ol	17	type	10		
20. ar	nalysis	17	analysis	11		
21. st	age	16	example	11		
22. tir	me	15				
23. re	action	15				
24. pr	oblem	15				
25. siı	mulation	14				

# Most frequent nouns in the 'this + noun' collocations

#### Comparing expert vs student writing: 'this + verb' collocations



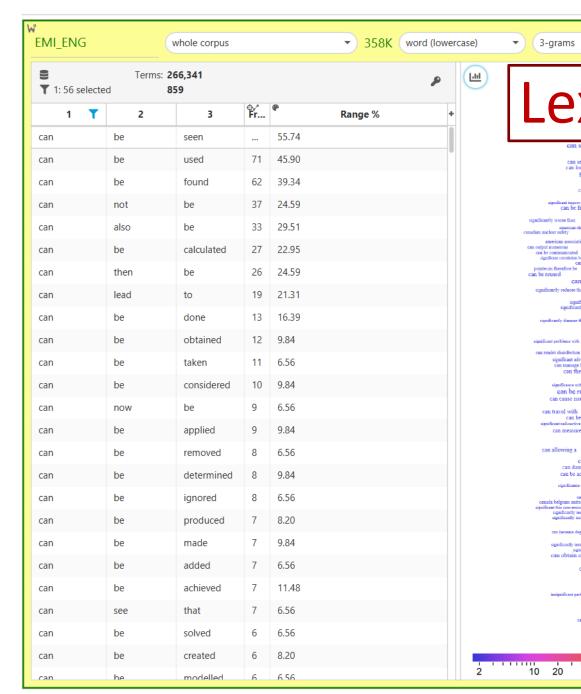
# Comparing more and less successful student writing

Linguistic category	Example	Relative frequency per 10,000 words		
	words	Lower-scored essays	Higher-scored essays	
Adjectival expressions of certainty	likely	156.7	196.9	
	possible	940.2	492.2	
	typical	52.2	69.2	
Indefinite pronouns for generalizability	anything	104.5	32.8	
	something	156.7	65.5	
	everyone	173.4	43.6	
Words for <b>hedging</b>	suggest	0	76.1	
	might	783.5	63.6	
Words for discourse	additionally	52.8	98.4	
organisation	thus	64.2	393.8	

## N-gram analysis (lexical bundles, bigrams,...)

What is the most common lexical bundle (3-word expression) starting with 'can' in the writing of students from Engineering?

CAN \_\_\_\_\_





□ ×



can be illustrated can be redimented can damage equipment significant levels of can become activated can just be can be better

1 1 1 1 1 1 1 1 1 1 1 1

## Exploring EMI corpus for EAP purposes



Contents lists available at ScienceDirect

#### Research Methods in Applied Linguistics

journal homepage: www.elsevier.com/locate/rmal



Gablasova, D., Harding, L., Bottini, R., Brezina, V., Ren, S., Savski, K., Iamartino, G., Li, Y., Liu, T., Poggesi, L., Toomaneejinda, A. & Zottola, A. (2024). Building a corpus of student academic writing in EMI contexts: Challenges in corpus design and data collection across international higher education settings. *Research Methods in Applied Linguistics*, 3(3), 100140.



# Challenges in data collection across diverse international education settings

## Discussion point

Please discuss the following question with a partner or a small group.

Have you been involved in data collection which involved multiple sites?

What were the benefits of the multi-site data collection?

What were the main challenges that you experienced?



## 1. A multi-site project: Different educational contexts

- A more comprehensive picture of the observed phenomenon
- Increased ecological validity of the findings compared to single-site research
- Findings from varied datasets can inform (pedagogical) practice across a wider variety of contexts
- Collaborative nature of such research and the ability to draw on the collective expertise of team members and their insights into local research sites

**Knowledge sharing** at different stages of the project:

- conceptualisation stage theoretical frameworks applicable to and inclusive of practices at different research sites;
- data collection enabling collaborators to share experience when issues arise,
- data analysis and interpretation the combined experience and expertise of team members can lead to "a more holistic understanding of findings" (Moranski & Ziegler, 2021, p. 223).

# A multi-site project: Challenges

Data collection logs to document challenges and strategies at each individual site

#### Problem (Aim)

Please use this section to describe and contextualize your issues/aims:

- 1. What was the aim you were trying to achieve?
- What were some key/different aspects of this aim?
- 3. What were the challenges encountered?

#### Solutions

Please use this section to record the strategies you used, and why they worked or did not work. You may address questions such as:

- What strategy/strategies have you used?
- How did the strategy/strategies work for you why did it work or didn't work?
- 3. What were the difficult aspects of solving the issue?
- 4. What helped you with dealing with this challenge?

# A multi-site project: Challenges

**Getting access** to research sites/participants – a potential challenge in any research with human participants - permissions required: **institutional level** & **level of different academic units** within the institution (e.g. faculty, department)

Multi-site research: permissions differed in scope and type across institutions involved in the project – difficult to anticipate/plan for

#### Example of requirements:

- In some cases, multiple levels of permission required within same institution e.g. at one research site, an approval was required from the faculty research unit, further approvals from various units within faculty, and an approval from the dean the same process was repeated for each faculty
- Different practices regarding ethical approval: some institutions accepted LU ethics, others required local ethical approvals

# A multi-site project: Challenges

Gaining access – required not only satisfying the administrative processes but also required permission from gatekeepers (eg. Deans, HoDs, teachers)

The procedure often not completely clear/straightforward

- the request for a permission could take a long time to be considered
- The permission depended not only on administrative procedures but also related to issues of trust, unfamiliarity with language-related research and perceived risks

**Impact on the project**: Issues with access affect ability to collect data in some disciplinary areas & some types of assignments

### Strategies for addressing challenges in gaining access

1. Being prepared to communicate the goals of the project to different audiences -> greater understanding of language-related research led to greater trust and cooperation

#### **Strategies**:

- written FAQ documents
- information/discussion sessions for staff in different departments
- recording short videos explaining the project
- showing examples of findings from corpus-based research
- showing examples of previous work completed by the researchers in the team

## Strategies for addressing challenges in gaining access

#### 2. Drawing on existing personal relationships:

- for gaining access to different institutional units (e.g. being able to come to a department to explain what we would like to do)
- shared contacts could help to 'vouchsafe' for the researchers/the project when
   establishing new contacts

#### 3. Prioritising personal, face-to-face communication:

- contacting students/departments via emails often led to delays
- personal, face-to-face meetings appeared more effective/efficient in long-term (helping to resolve issues of trust, familiarity with linguistics research, etc)

### 2. Challenges 'on the ground': Recruiting students

#### Two major challenges:

- Establishing initial contact
- Gaining consent and obtaining the data

#### **Establishing contact** with students & **explaining the project**:

- the need for different context-appropriate strategies
- the strategies differed according to the country, institution, academic unit
- required flexibility and creativity

# Challenges 'on the ground': Recruiting students

**Strategies**: contacting students via departments, using financial incentives (ranging from Amazon vouchers, honoraria, book tokens, coupons for coffee/McDonalds/KFC breakfasts/movies, price draws), contacting students via student reps, social groups; organising information sessions about the project, recording videos and sharing them with students.

While multi-site design made this more challenging – it was also a great source of **solutions** - the combined expertise of the team and understanding of the local context were crucial

- Good understanding of local culture and values
- Sharing ideas about strategies

## Discussion point

#### Please discuss the following questions with a partner/small group:

- As part of your research, did you need to obtain ethics permission for collecting data outside of your home country/institution? If yes, how long did it take to obtain the required permissions?
- Did you experience any challenges when accessing data/participants outside of your home institutions – how did you approach these?



# 3. Construct of student academic writing

- Decisions about what language samples to include in a corpus are central in corpus design → implications for representativeness and generalizability
- Aim of current project compile a corpus of student writing from different universities and countries → we needed a construct of academic writing that can be meaningfully applied across different higher education institutions



Prince of Songkla University



Thammasat University



University of Turin



University of Milan



Xi'an Jiaotong University



Xi'an Jiaotong-Liverpool University

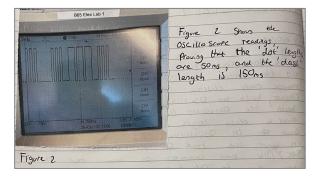


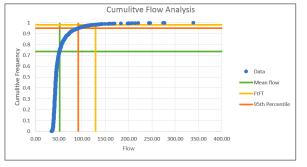
Lancaster University

#### Operationalising student academic writing: Challenges

- Academic writing is a complex notion can refer to and encompass very varied set of writing practices – related to the enormous diversity of academic actors, communicative aims, values, motivations, etc in academic study and research (Hyland, 2006)
- Student writing: formal assessed pieces informal notes written during group discussions – emails to course tutors – lecture notes - etc
- Specific operationalisation of the construct → impact on the selection/inclusion of texts → impact on the type of academic writing represented (or excluded) in the corpus

An intense rainfall, carthquake shaking, volcanic eruption, storm waves, or rapid stream
crosion are causes of increasing the stresses and reducing the strength of slope materials which
triggers landslides (Wieczorek 1996). It is also anticipated that incidents of land slide disasters
may possibly increase due to over exploitation of natural resources, rapid deforestation, climate
change, and increase in hill population and uncontrolled excavations which results in higher
susceptibility of surface soil to instability (Manivannan and V. Kasthuri 2020). Van et al. (2010,
also add that it is assumed that natural factor are considered as prime factor for the landslide an
human activities are considered as less important. Human are regarded as victim of landslide an
are considered vulnerable to the disaster but not studied as a factor that might be responsible for





### EMI Corpus: Construct of student academic writing

- All institutions: Disciplinary writing, submitted for assessment
- Institutions differed in their preferred types of writing: Electronic vs handwritten
  - Differences in the type of writing practices and processes (e.g., editing, planning, access to resources, exam setting, effect of stress)

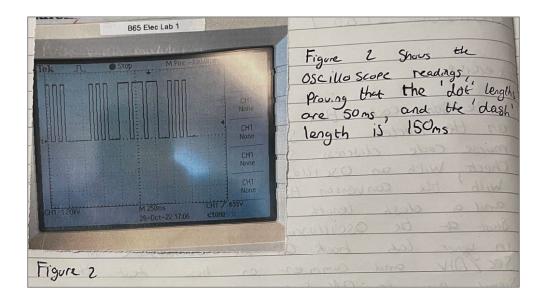
#### Construct of academic writing: Different writing practices

An intense rainfall, earthquake shaking, volcanic eruption, storm waves, or rapid stream
erosion are causes of increasing the stresses and reducing the strength of slope materials which
triggers landslides (Wieczorek 1996). It is also anticipated that incidents of land slide disasters
may possibly increase due to over exploitation of natural resources, rapid deforestation, climate
change, and increase in hill population and uncontrolled excavations which results in higher
susceptibility of surface soil to instability (Manivannan and V. Kasthuri 2020). Van et al. (2010)
also add that it is assumed that natural factor are considered as prime factor for the landslide and
human activities are considered as less important. Human are regarded as victim of landslide and
are considered vulnerable to the disaster but not studied as a factor that might be responsible for

Innovation is the process of creating a se product, process, oraganizational methods on marketing methods in a new on improved way of that existing identity product, process, oraganizational methods on marketing

Writing practices reflecting different contexts of production → typical linguistic features

- Handwritten vs electronically submitted
- Produced in timed vs non-timed conditions
- Produced under exam conditions

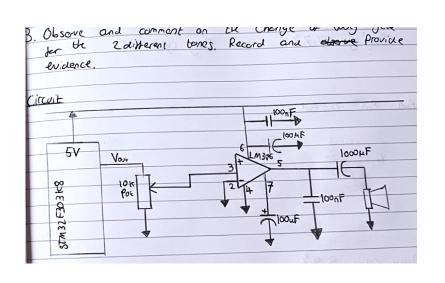


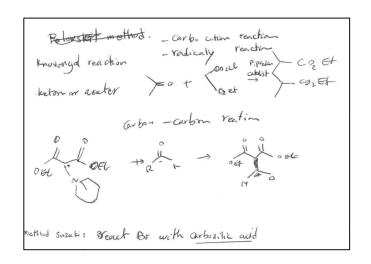
### EMI Corpus: Construct of student academic writing

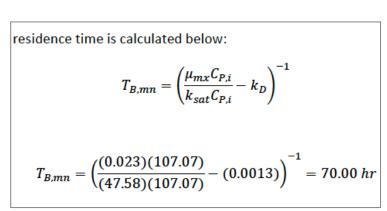
- Disciplinary writing, submitted for assessment
- Electronic & handwritten submissions
- Word length varied considerably across disciplines and institutions: Written pieces
  - min. 100 words including text, figures, diagrams, code, etc.

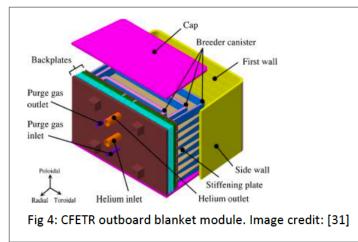
#### Construct of academic writing: Different writing practices

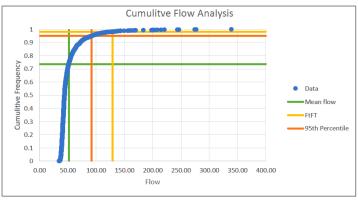
Capturing the **visual aspects** of student production → insights into the changing nature of what counts as 'academic writing' and in what way this differs across disciplines (e.g. STEM subjects)











## EMI Corpus: Construct of student academic writing

- Disciplinary writing, submitted for assessment
- Electronic & handwritten submissions
- Written pieces min. 100 words including text, figures, diagrams, code

## Construct of student academic writing

The adopted construct prioritised – as much as possible – an **inclusive approach** – to maximise the opportunities offered by access to multiple educational sites

The final approach followed close discussions with partner universities to understand:

- i) the type of writing produced but also ii) the status of different types of writing and
- iii) the local writing practice/needs

## Discussion point

Addressing challenges related to deciding on the construct of 'academic writing' in research/teaching

Please discuss the following questions with a partner or a small group:

- What type of student academic writing would you like to see represented in research/corpora?
- Are there particular challenges when collecting the target data across different institutions?

# Concluding thoughts

The session highlighted...

- some of the challenges involved in a multi-site, international project involving corpus construction and the strategies/ approaches used to address them
- the interaction of theoretical, methodological and practical considerations that are part of collaborative projects that involve data collection at different sites
- the value of (large) collaborative projects that can draw on the experience and expertise of researcher from different educational backgrounds and settings

## Thank you!

#### Corpus research informing EMI practice





https://wp.lancs.ac.uk/emi-corpus-project/