



# Harmony

Open source NLP/AI tool for psychologists and social and health sciences  
Data discovery and harmonisation

[harmonydata.ac.uk](http://harmonydata.ac.uk)  
[github.com/harmonydata/](https://github.com/harmonydata/)

Thomas Wood  
 **Fast Data Science**

Bettina Moltrecht, PhD  
Centre for Longitudinal Studies  
UCL



Economic  
and Social  
Research Council



Economic  
and Social  
Research Council

# MEET THE HARMONY TEAM



**Rachel Gomes**  
University College London



**Richard Thomas**  
UK-LLC



**Eoin McElroy**  
Ulster University



**Bettina Moltrecht**  
University College London



**Louise Arseneault**  
Kings College London



**Thomas Wood**  
Fast Data Science



**Mauricio Hoffmann**  
Universidade Federal de Santa Maria

CENTRE FOR  
LONGITUDINAL  
STUDIES

Ulster University

UK  
Longitudinal  
Linkage  
Collaboration

Fast Data Science





# CONTRIBUTORS

...plus people around  
the world who have  
made 24 pull requests  
to Github!

**vkrithika25  
makrianast  
shahid-0  
EveWCheng  
0x48piraj  
nikhildevre  
omtarful  
Gabe-ferr  
abdullahwaqar  
deannavarley**

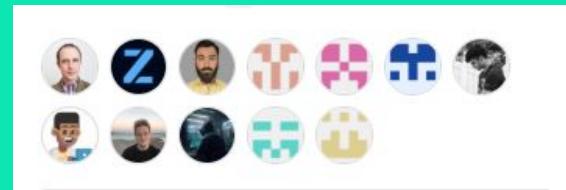
Filters ▾ Q is:pr is:closed

Clear current search query, filters, and sorts

2 Open ✓ 19 Closed

- Allow batching of items when sent to LLM ✓ #66 by makrianast was merged 4 days ago 10 tasks
- #62 Added crosswalk table + unit tests ✓ #65 by vkrithika25 was merged 5 days ago 10 tasks done
- feat(utils): Add strip\_prefixes function to remove common question words ✓ #64 by abdullahwaqar was merged 4 days ago 10 of 13 tasks
- Search instruments ✓ #52 by zaironjacobs was merged on Sep 6 1 task done
- Upgrade to Pydantic 2.8.2 ✓ #51 by olp-cs was merged on Aug 20 6 of 11 tasks
- Replace PDF parsing ✘ #49 by woodthom2 was merged on Jul 19 12 tasks done
- Catalogue match instruments ✓ #48 by zaironjacobs was merged on Aug 21
- Load instrument from list ✓ #47 by woodthom2 was merged on Jun 21 3 tasks

SECTION



# Partners



Economic  
and Social  
Research Council









The challenge...



# Combining studies...

## Retrospective harmonisation

We have to match different survey items, sometimes in different languages

**GAD-7 Anxiety**

Over the <u>last two weeks</u> , how often have you been bothered by the following problems?	Not at all	Several days	More than half the days
1. Feeling nervous, anxious, or on edge	0	1	2
2. Not being able to stop or control worrying	0	1	2
3. Worrying too much about different things	0	1	2
<b>4. Trouble relaxing</b>	0	1	2
5. Being so restless that it is hard to sit still	0	1	2
6. Becoming easily annoyed or irritable	0	1	2
7. Feeling afraid, as if something awful might happen	0	1	2

Column totals \_\_\_\_\_ + \_\_\_\_\_ + \_\_\_\_\_  
Total score \_\_\_\_\_

**Beck's Anxiety Inventory**

- Numbness or tingling
- Feeling hot
- Wobbliness in legs
- Unable to relax**
- Fear of worst happening
- Dizzy or lightheaded
- Heart pounding / racing
- Unsteady
- Terrified or afraid
- Nervous
- Feeling of choking
- Hands trembling
- Shaky / unsteady
- Fear of losing control
- Difficulty in breathing
- Fear of dying
- Scared
- Indigestion
- Faint / lightheaded
- Face flushed
- Hot / cold sweats



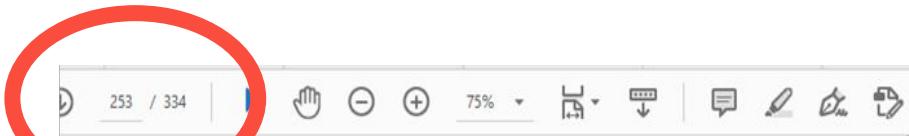
# An example...

Research using data from the UK and Brazil

UK: Millennium cohort study  
19,000 children

Brazilian High Risk Cohort Study for Childhood Mental Health Conditions (BHRCS).  
2,511 young people

Surveys were in different languages, using different instruments



A1	B	C
<b>BHRCS Parent report ABCL</b>		
Pergunta		
ABCL5a	Aproximadamente quantos(as) amigos(as) próximos(as) ele(a) tem? (Não incluir pessoas da família)	Nível de resposta
ABCL5b	Aproximadamente, quantas vezes por mês ele(a) tem contato com qualquer um dos amigo(a) próximo(a)? (Incluir contatos pessoalmente, por telefonemas, cartas, e-mail)	0-99
ABCL5c	Até que ponto ele(a) se dá bem com amigos(as) próximos?	0-99
ABCL5d	Aproximadamente, quantas vezes por mês amigos ou familiares o (a) visitam?	---
ABCL6	Qual o estado civil dele(a)?	
ABCL7	Em algum momento nos últimos seis meses, ele(a) viveu com um(a) esposo(a) ou con-	
ABCL7a	Agora por favor, circule 0, 1 ou 2 ao lado das afirmações para descrever o rela-	
ABCL7a1	Se dá bem com o esposo(a) ou companheiro(a)	1
ABCL7a2	Tem dificuldade em dividir responsabilidades com esposo(a) ou companheiro(a)	2
ABCL7a3	Parece satisfeito(a) com esposo(a) ou companheiro(a)	3
ABCL7a4	Gosta das mesmas atividades que o esposo(a) ou companheiro(a)	4
ABCL7a5	Discorda do esposo(a) ou companheiro(a) sobre questões de administração do lar, co-	All
ABCL7a6	Ele(a) tem problemas com a família do esposo(a) ou companheiro(a)	All
ABCL7a7	Gosta dos amigos do esposo(a) ou companheiro(a)	1 9 months
ABCL7a8	Incomoda-se com o comportamento do esposo(a) ou companheiro(a)	2 3
ABCL8	Ele(a) tem alguma doença ou deficiência?	1 9 months
ABCL8a	Se sim, por favor, descreva:	Ethnicity
ABCL9	Por favor, descreva qualquer preocupação que você tenha sobre ele:	2 3
ABCL10	Por favor, descreva as principais qualidades e pontos positivos dele(a):	1 9 months
24		Social class at birth
25	[Entrevistador] Mostrar C17	Social class at birth
26	[LEIA] Logo abaixo, você encontrará uma lista de afirmações que descrevem as pessoas. Para cada afir-	Social class at birth
27		Social class at birth
28	Número	Pergunta
29	AL1	É muito esquecido(a)
30	AL2	Sabe aproveitar as oportunidades que lhe aparecem
31	AL3	Discute muito
32	AL4	Desenvolve suas habilidades
33	AL5	Culpa os outros por seus próprios problemas
34	AL6	Usa drogas (que não álcool ou nicotina) sem fins medicinais (descreva):
35	AL7	Gosta de contar vantagem
36	AL8	Não consegue concentrar-se, não consegue prestar atenção muito tempo
37	AL9	Não consegue tirar certos pensamentos da cabeça; obsessões (descreva):
38	AL10	Não consegue parar sentado(a), não parar quieto(a) ou é hiperativo(a)

+ BHRCS Parent report ABCL ▾ BHRCS\_Parent

24	1 9 months	Breastfeeding	ambfeda (ambfedaa0 ambfedab0 ambfedc0 in original wide format)
25	1 9 months	Breastfeeding	ambfed (ambfedaa0 ambfedab0 ambfedc0 in original wide format)
26	1 9 months	Breastfeeding	ambfew (ambfewaa0 ambfewab0 ambfewc0 in original wide format)
27	1 9 months	Breastfeeding	ambfem (ambfema0 ambfemb0 ambfemc0 in original wide format)
28	2 3	Breastfeeding	bmbfmt (bmbfma0 bmbfmb0 bmbfmc0 in original wide format)
29	2 3	Breastfeeding	hmhfrea (hmhfreaa0 hmhfreh0 hmhfrc0 in original wide format)

+ Sheet1 ▾

253 / 334



## MAIN AND PARTNER RESPONDENTS

### INWB

Here is a question about your feelings on how satisfied you are with your life. There is no right or wrong answer. Please give an answer on the scale from 0 to 10 where '0' means completely dissatisfied and '10' means completely satisfied.

SATN Overall how satisfied are you with your life nowadays?

### SCALE FOR WELLBEING GRID:

'Scale: 0-10'

0 completely dissatisfied

10 Completely satisfied

11 Don't know/don't wish to answer

age | 252

### PAUSE

### SCEN

Thank you for answering these questions. Please tell the interviewer you have finished now and they will carry on with the interview.

### PAUSE

PLEASE RE-ATTACH THE TABLET TO THE KEYBOARD. PLEASE BE CAREFUL NOT TO TOUCH THE SCREEN WHEN YOU DO THIS.

### SCFI

INTERVIEWER CODE (DO NOT ASK). DID THE RESPONDENT ANSWER ALL OF THE QUESTIONS IN THIS SECTION VIA CASI SELF-COMPLETION?

1 Yes, all self-completion by respondent

2 Yes, self-completion, but interviewer helped to complete some questions

3 No, interviewer completed it all with the respondent.

[Don't Know and Refusal are not allowed]

ND OF FILTER

A1	Cohort																	
1	Cohort	Variable name	wave	Question	Question in Portuguese	Response option	Recoding	Comments	Name (harmonized)	Final coding	Recoding/Comments	Response option	Question	wave	Variable name	Comments		
PARENT-CHILD RELATIONSHIP (in contact: yes=0/no=1; quality: 3= not close, 2=somewhat, 1= very close 0= extremely close)																		
10	BHRCs	mrelation (reg22) w0	dom_s01_reg_w0	How does/did the child and the biological mother get along when they are/were together?	Como a criança e mãe biológica se dão ou se davam quando estão/estavam juntas?	1, Very well 2, Well 3, More or less 4, bad 5, Very bad 88, Does not apply 99, Does not know	Brazil 4 5 = 3 Brazil 3 = 2 Brazil 2 = 1 Brazil 1 = 0	mrelation	In MCS this is mostly mothers, but can also be father/other caregiver, so we either combine Brazil or do only	3= not close, 2=somewhat close, 1= very close 0= extremely close	UK 4 = 0 UK 3 = 1 UK 2 = 2 UK 1 = 3	1 Not very close 2 Fairly close 3 Very close 4 Extremely close 5 Dont know/Dont wish to answer	Overall, how close would you say you are to CM? (7)Overall, how close would you say you are to your father?	5, 6, 7	EPSCHC00 FPSCHC00 (parent_cm interview) FCRLQM00 (cm interview) FCRLQF00 (cm interview) GCRLQF00 (cm interview) GPSCHC00parent_cm interview)			
Marital status																		
14	BHRCs	reg31	dom_s01_reg_w0w2	What is the marital status of the biological mother?	A mãe biológica da criança está no momento: casada ou morando junto com o pai biológico, casada ou morando junto com companheiro,	88, Refuses 99, Does not know 1, Married or living together with the biological	Brazil 1= 1 Brazil 2= 2 Brazil 3 4= 3 Brazil 5= 4	maritalstatus	1= Married or living with biological parent 2= Separated and living with someone else 3= Separated /Divorced	I think the recoding needs changed as we do not have Separated and living with someone else easily available in MCS	Refused   -9 Dont know   -8 Not applicable   -1 Legally separated   1 Married, 1st and only marriage   2 Remarried, 2nd	Current legal marital status	5, 6, 7	EPFCIN00 (parent interview) FPFCIN00 (parent interview) GPFCIN00 (parent interview)				
PARENT DEATH (yes=1, no=0)																		
16	BHRCs	reg36	dom_s01_reg_w0w2	Is the child's biological father known?	O pai biológico da criança é desconhecido, conhecido ou falecido?	1, Unknown 2, Known 3, Deceased	Brazil 3 = 1, Brazil 1 2 = 0	fatherstatus	0= not dead 1= dead	Coding options reduced as of wave 6	3=1 "dead" all other =0	1 Resident full-time in household   2 Resident part-time in household   3 Deceased   4 Non-resident, in contact   5 Non-resident, not in contact   6 Non resident, contact	Natural father status	5,6,7	ENATFOO Family derived FDNATFOO (family derived) GDNATFOO (family derived)			



... months later ...



# HARMONY

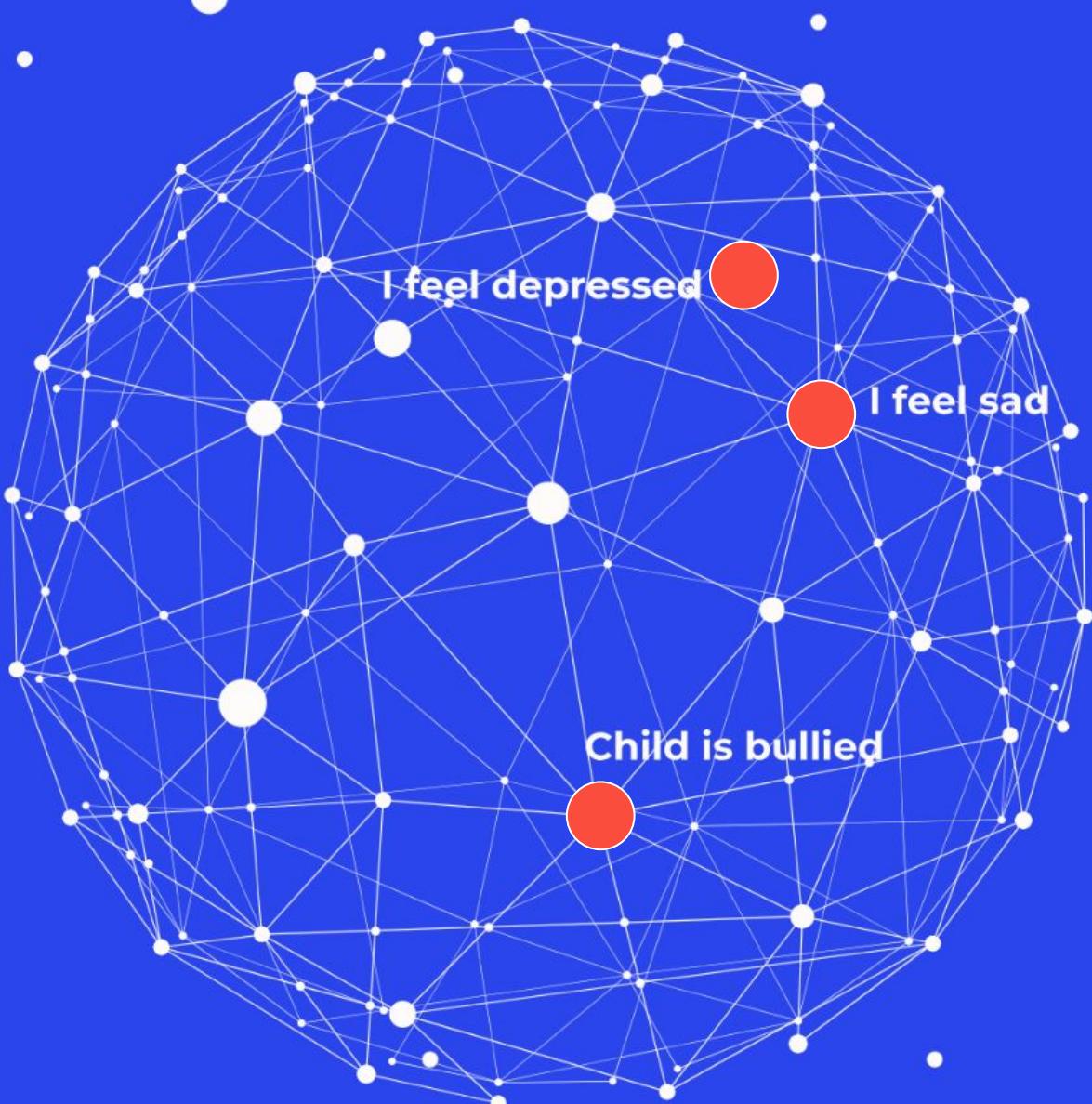
<https://harmonydata.ac.uk/app>



# How does Harmony work?

1. Process the PDF to get the question items out
2. Convert all questions to sentence embeddings
3. Calculate the cosine similarity

<https://fastdatascience.com/natural-language-processing/sematic-similarity-with-sentence-embeddings/>



Sentences have a related meaning

Vectors point in a similar direction

Cosine score is high – close to 1

Sentences have unrelated meaning

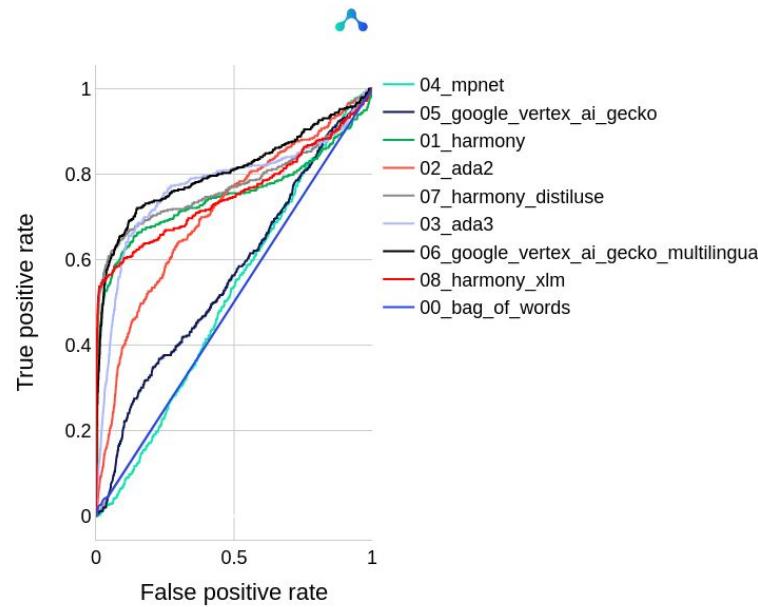
Vectors are orthogonal

Cosine score is low – closer to 0

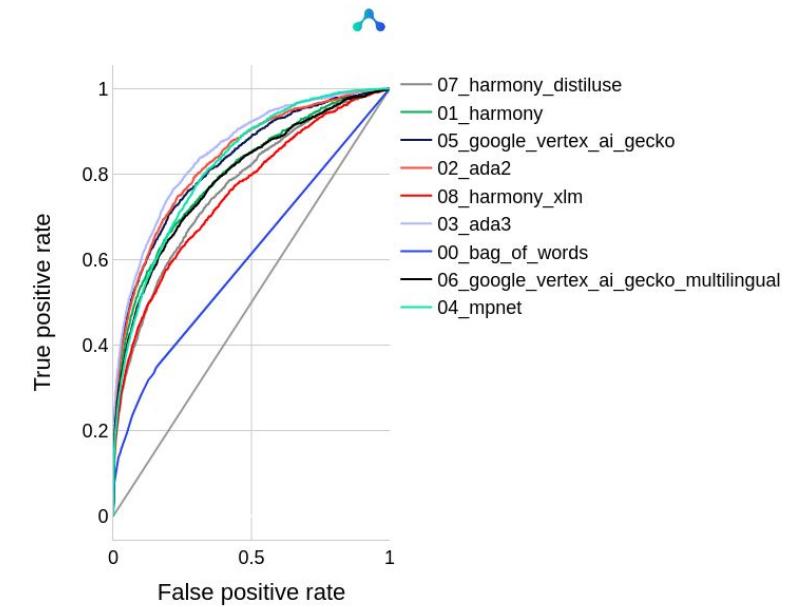


# Evaluating Harmony

ROC on GAD 7 multilingual dataset



ROC on McElroy et al childhood dataset



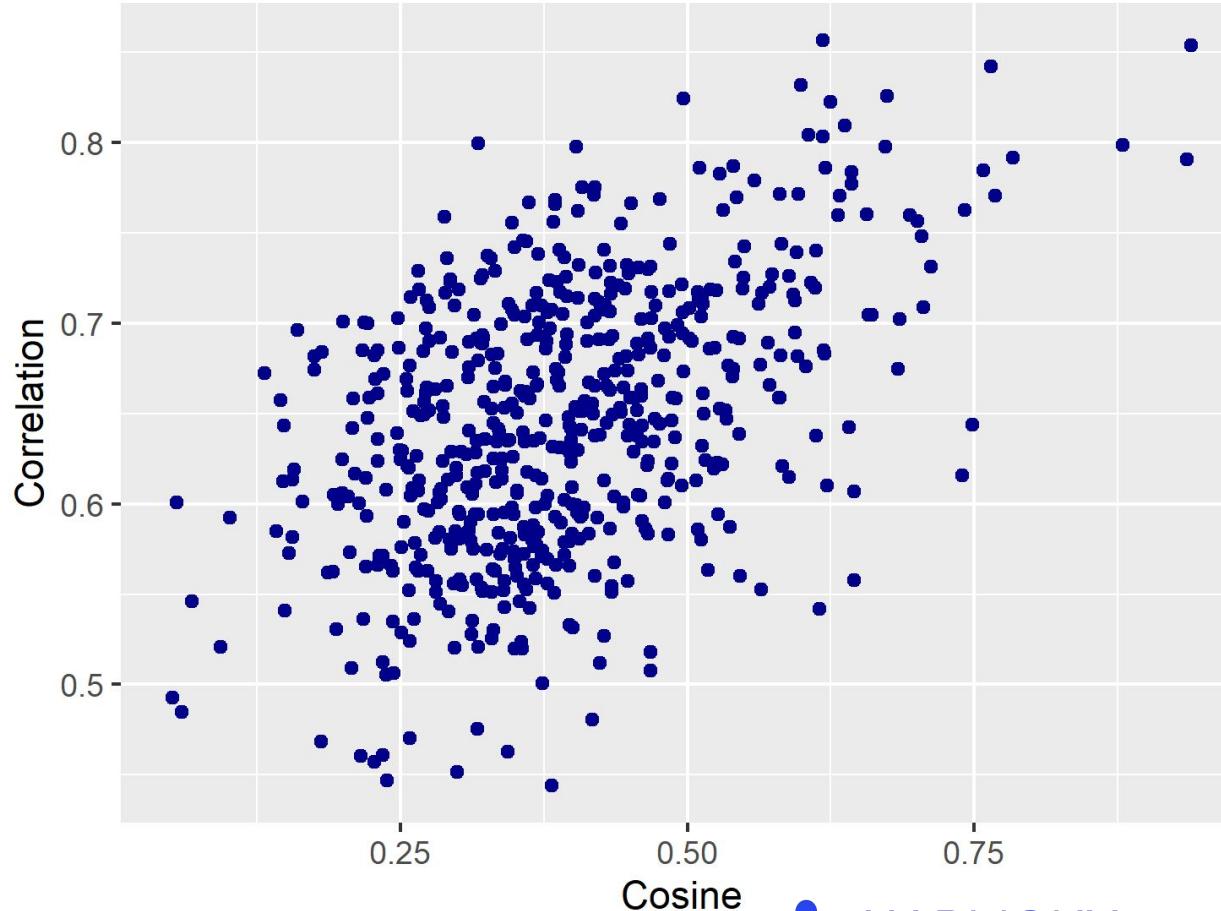


# Real correlations?

	A	B	C					
1	Questionnaire	Item number	Content					
2	IDQ							
3	IDQ							
4	IDQ							
5	IDQ							
6	IDQ							
7	IDQ							
8	IDQ							
9	IDQ							
10	IDQ							
11	IAQ							
12	IAQ							
13	IAQ							
14	IAQ							
15	IAQ							
16	IAQ							
17	IAQ							
18	IAQ							
19	DHQ							
SECTION								
			A		B	C	D	E
			1 Supplementary File 2. Correlaiton and cosine coefficients for item pairs		from	to	spearman	cosine
					1	10	0.719538559	0.61149627
					2	11	0.719244021	0.445720732
					3	12	0.731182941	0.711875081
					4	13	0.665979411	0.571581244
					5	14	0.703580795	0.511955619
					6	15	0.709961188	0.297138691
					7	16	0.691486691	0.50184983
					8	17	0.608636306	0.259960353
					9	18	0.738485659	0.37026453
					10	19	0.798613012	0.879561961
					11	2	0.824578169	0.496798843
					12	20	0.573991369	0.572757006
					13	21	0.628075543	0.655307412
					14	22	0.625625193	0.354875147
					15	23	0.718334673	0.681931853
					16	24	0.683265023	0.334293664
					17	25	0.587003714	0.53709048



# Real correlations?



MCElroy et al,

Using natural language processing to facilitate the harmonization of mental health questionnaires: a validation study using real-world data



# Harmony Discovery

<https://harmonydiscovery.fastdatascience.com/>

We are allowing researchers to discover datasets and studies by searching in a vector index

A giant vector index of UK studies across disciplines  
Currently using Elasticsearch (experimenting with Pinecone)

Designs: [https://www.figma.com/...](https://www.figma.com/)



# Open source

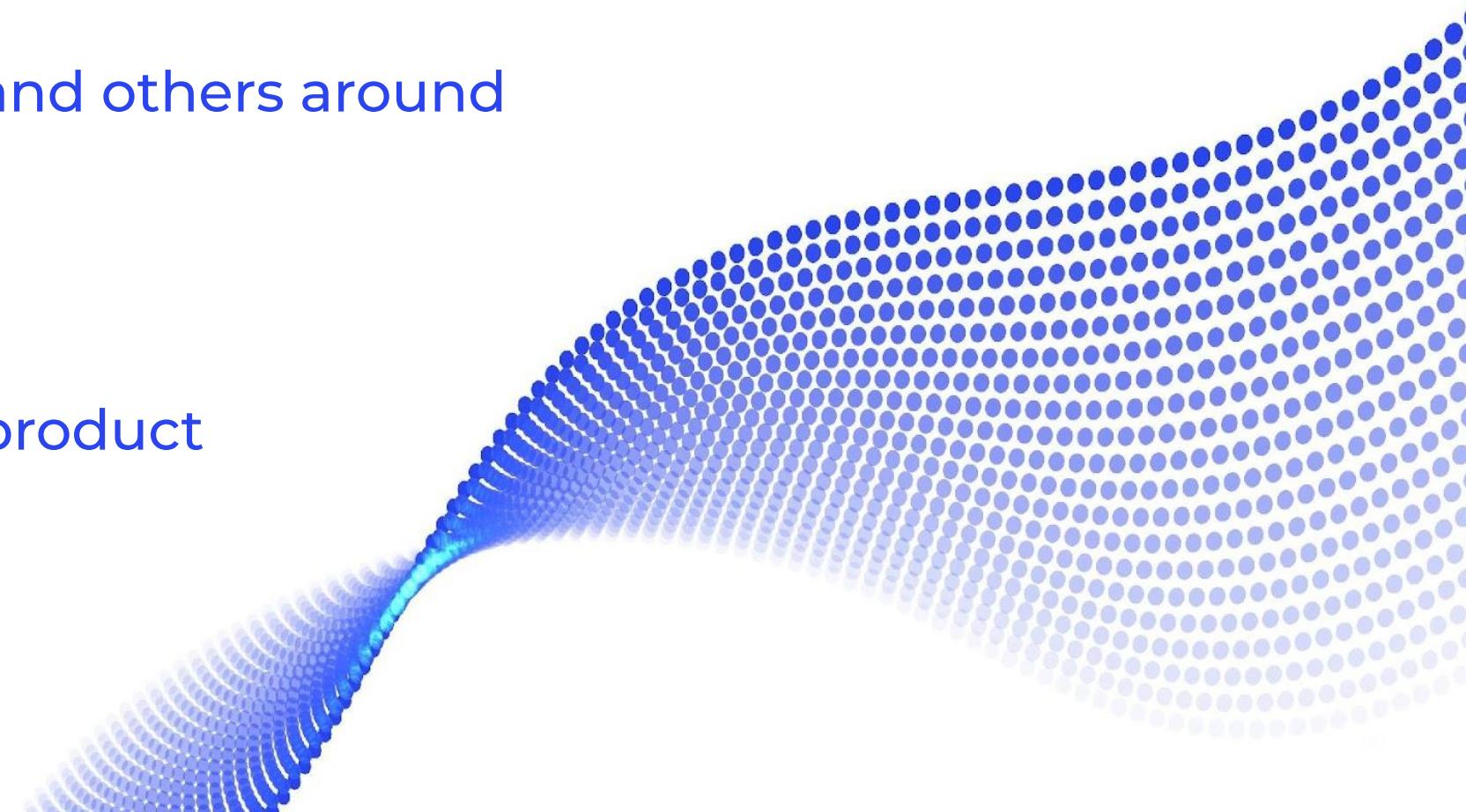
<https://github.com/harmonydata>

Free for researchers and others around  
the world

MIT License

It's not a monetised product

(online) hackathons





# Discord

discord.gg/harmonydata

general 6 October 2024

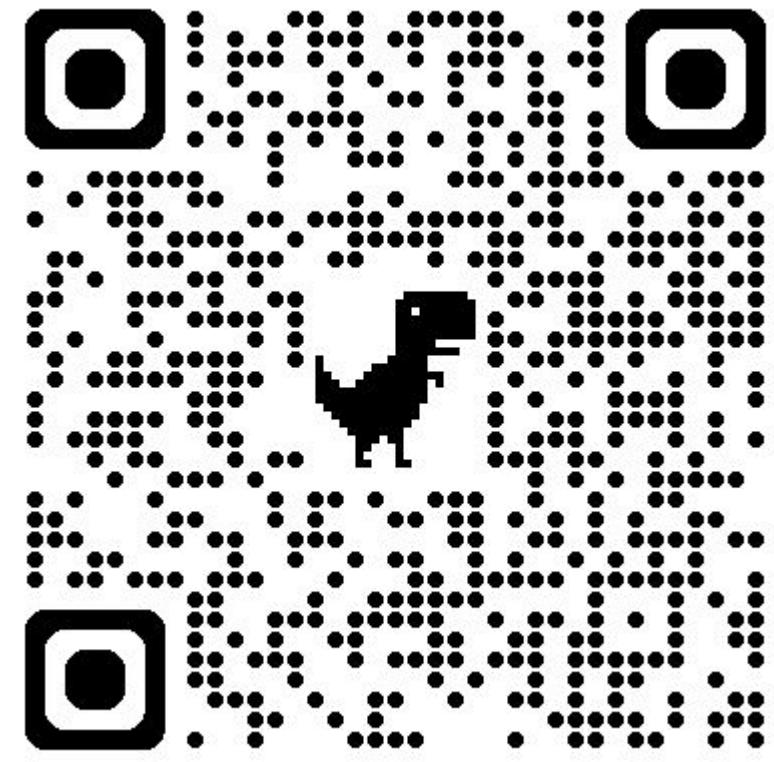
**Lucious** 03/10/2024 13:22  
Hello, i am from türkiye, i am 20 years old,i am a computer engineering student and i am still studying. i learned about Harmony when i met a person and mentioned it to me. (edited)  
and the project attracted my attention, I want to be useful

**Thomas** 06/10/2024 09:54  
Hi Lucious  
thanks, I'm glad to have you on board!  
We have some ideas for what you could do if you would like to pick up a task:  
1. PDF report export function – can we make Harmony create a shareable PDF export? – you could use code from this open source app  
<https://app.clinicaltrialrisk.org/>

**Clinical Trial Risk Tool**  
Analyse your Clinical Trial protocols and identify risk factors using Natural Language Processing, from Fast Data Science.

2. Ability to harmonise data (can we create a Notebook on the fly?) – can you make it so that the user can export Harmony's analyses to a Colab or Jupyter notebook so that they can then analyse the data?  
do you have Python ability or back end/front end skills?  
Also **@everyone** On Tuesday next week, I'm giving a presentation at an AI meetup called AI|DL in London - if you're in the UK you might be interested to come.  
Here's the info about the meetup: <https://harmonystatistics.ac.uk/psychology/aidl-meetup/>

Startup  
**Harmony at AI|DL meetup | Harmony**  
A global platform for contextual data harmonisation

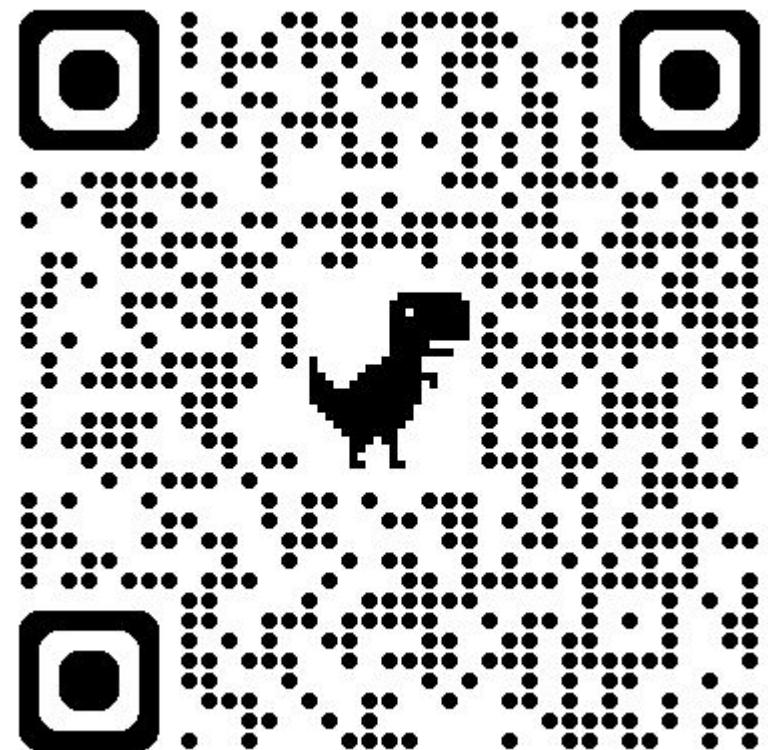




# Open source tasks

We are always looking for collaborators.  
Ideas include

- Create browser plugin
- Integration with more platforms
- PDF report export function (can we make Harmony create a shareable PDF export?)
- Ability to harmonise data (can we create a Notebook on the fly?)



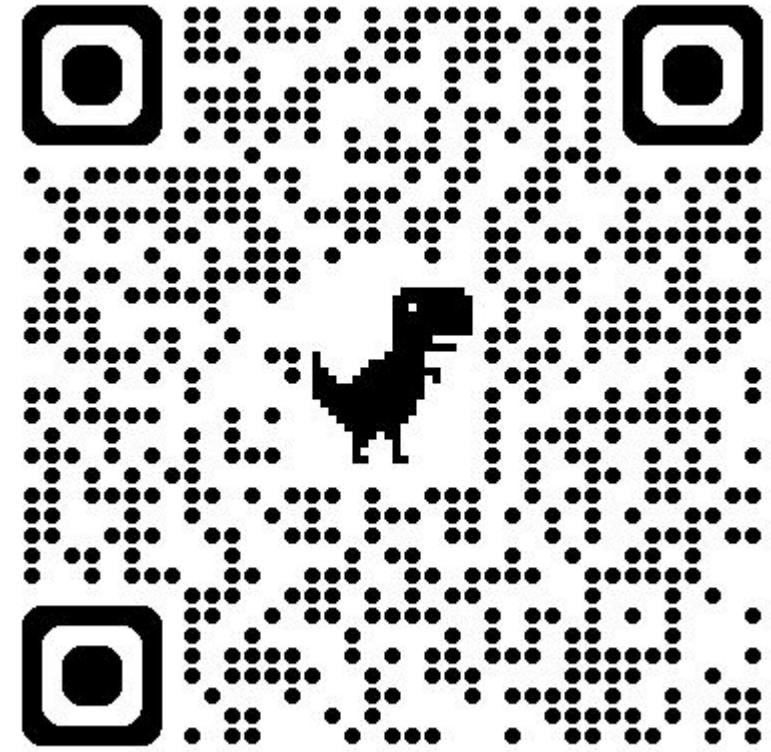
Search

#	Participant	Submitted	MAE	R <sup>2</sup>
1	<u>Alex</u> @ayetsu	11 hours ago	20.683	0.14331
2	<u>Daksh</u> @Daksh	1 day ago	20.906	0.16552
3	<u>Akshat</u> @akshat	9 hours ago	21.198	0.08753
4	<u>Jeremy.</u> @jezz	1 month ago	21.261	0.11211
5	<u>Gabriele Monti</u> @montig	2 weeks ago	21.629	0.04189
6	<u>idle</u> @Idlehunter	2 weeks ago	21.770	0.01900
7	<u>Jake</u> @jakedorman64	1 month ago	21.797	0.09679
8	<u>Louis de Wardt</u> @louisdewardt	3 weeks ago	21.797	0.09679
9	<u>Denis</u> @denis_chaykovskiy	2 weeks ago	21.797	0.09679
10	<u>Rafi Ahmed Patel</u> @rafa	6 days ago	21.797	0.09679
11	<u>Thomas</u> @thomas	1 month ago	24.666	-0.62098



# Train PDF parser

The screenshot shows the DOXA AI competition interface. At the top, there's a navigation bar with 'DOXA AI' and links for 'Home', 'Competitions', and 'Submissions'. A blue header banner displays the title 'Questionnaire PDF Parsing Challenge' and a subtext: 'A competition to help researchers extract and tag survey questions from PDFs accurately'. Below the banner, it says 'Harmony · Starting on 20 Jan 2025, 10:00:00 GMT'. A teal 'Overview' button is visible. The main content area features the title 'Harmony Questionnaire Parsing Challenge' with the HARMONY logo. To the right, there are competition details: 'Competition tag: harmony-parsing', 'Submission size limit: 4,295 MB', and a '23 members online' button.





# Thank you for listening!

 [harmonydata.ac.uk](http://harmonydata.ac.uk)

 [discord.gg/harmonydata](https://discord.gg/harmonydata)

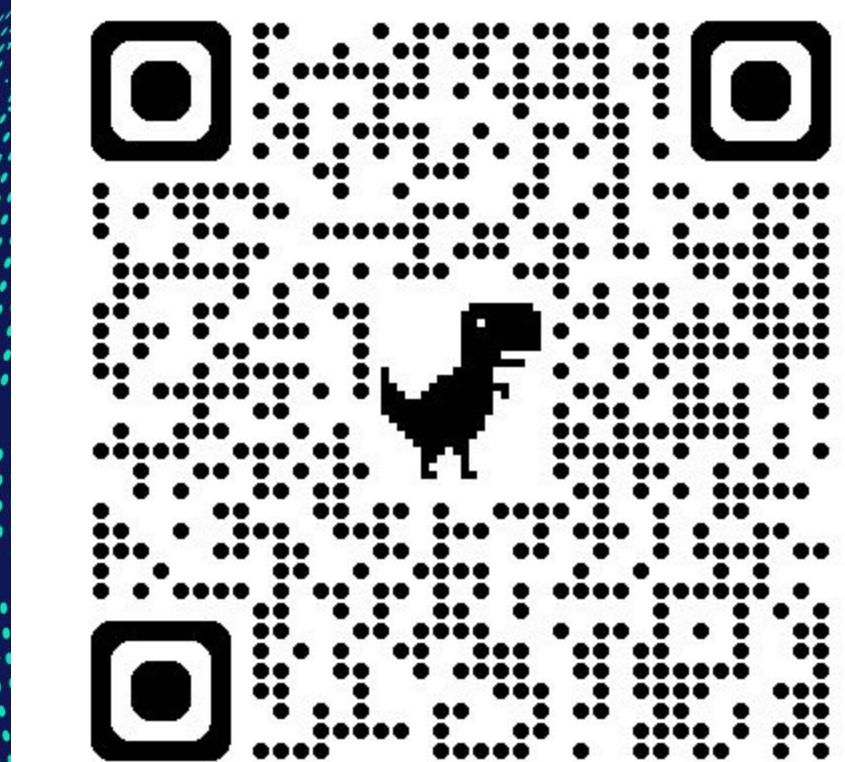
 [linkedin.com/company/harmonydata](https://linkedin.com/company/harmonydata)

 [github.com/harmonydata/harmony](https://github.com/harmonydata/harmony)

 [@harmony\\_data](https://twitter.com/@harmony_data)

 [@harmony\\_data](https://mastodon.social/@harmony_data)

[b.moltrecht@ucl.ac.uk](mailto:b.moltrecht@ucl.ac.uk) [thomas@fastdatascience.com](mailto:thomas@fastdatascience.com)





# HARMONY

## Harmonise questionnaire items - with Harmony

Harmony is a tool for retrospective harmonisation of questionnaire items.

If you want to compare items from different surveys, such as GAD-7 and PHQ-9, Harmony can identify which questions match.

[FAQs](#) - [Privacy policy](#)



Upload or drag and drop any **pdf, docx or xlsx** file here

Choose from existing instruments

SCARED English, GAD-7 English, CES\_D English

▼ SCARED English



▼ GAD-7 English



▼ CES\_D English



**HARMONISE**



## Options

Match Threshold



Search

Show within-instrument matches



SHARE

EXPORT

Found 8 matches

SCARED English - Q28

GAD-7 English - Q8.0

People tell me that I worry too much

81

Worrying too much about different things

GAD-7 English - Q7.0

CES\_D English - Q10

Feeling afraid, as if something awful might happen

78

I felt fearful.

SCARED English - Q85

GAD-7 English - Q3.0

I worry about how well I do things

76

Worrying too much about different things

SCARED English - Q15

GAD-7 English - Q7.0

When I get frightened, I feel like things are not real

74

Feeling afraid, as if something awful might happen

SCARED English - Q6

GAD-7 English - Q7.0

When I get frightened, I feel like passing out

72

Feeling afraid, as if something awful might happen

SCARED English - Q9

GAD-7 English - Q1.0

People tell me that I look nervous

71

Feeling nervous, anxious, or on edge

SCARED English - Q12

GAD-7 English - Q7.0



# HARMONY

## Options

Match Threshold



Search

fam

Show within-instrument matches



SHARE

EXPORT

Found 113 matches

SCARED English - Q29

I don't like to be away from my family

46

CES\_D English - Q8

I felt that I could not shake off the blues even with help from my family or friends.

SCARED English - Q29

I don't like to be away from my family

39

CES\_D English - Q20

I could not get "going."

SCARED English - Q29

I don't like to be away from my family

38

CES\_D English - Q19

I felt that people dislike me.

SCARED English - Q29

I don't like to be away from my family

37

CES\_D English - Q1

I was bothered by things that usually don't bother me.

GAD-7 English - Q2.0

Not being able to stop or control worrying

36

CES\_D English - Q8

I felt that I could not shake off the blues even with help from my family or friends.

SCARED English - Q29

My child doesn't like to be away from his/her family

34

CES\_D English - Q8

I felt that I could not shake off the blues even with help from my family or friends.

SCARED English - Q81

CES\_D English - Q8