Allied Data Science Communities

# Allied Data Science Communities

DON'T PANIC, WE ARE DATA SCIENTISTS!

# KAGGLE COVID-19 CHALLENGE INTRO

FRANKFURT DATA SCIENCE

VIENNA DATA SCIENCE GROUP

BUDAPEST.AI

BUDAPEST DEEP LEARNING READING SEMINAR

# LEVENTE SZABADOS

"...originally Buddhist theologian and programmer, AI professional (NLP), lead of research, lecturer,  startupper, ex-CTO
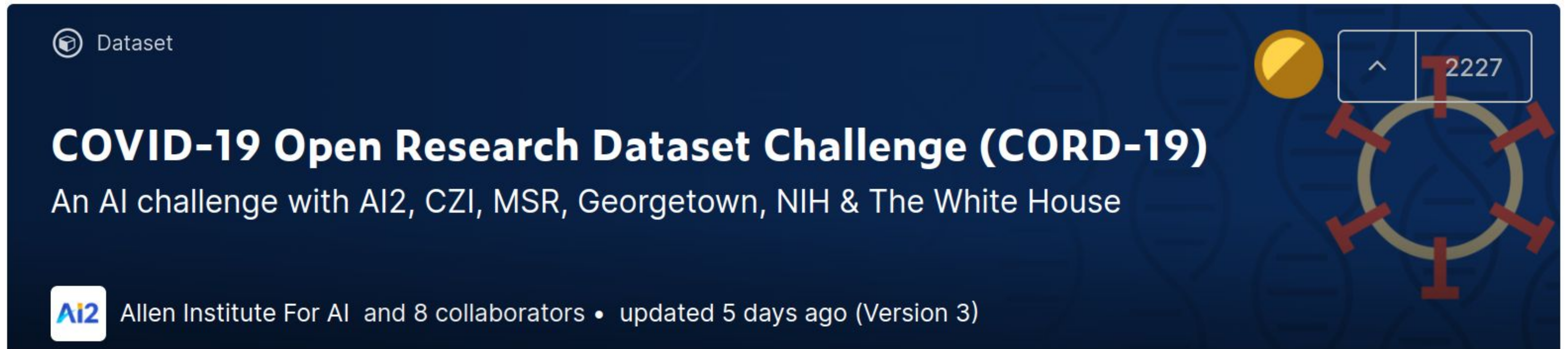
Presently:

Lecturer: Frankfurt School of Finance and Management,

Specialization leader: KÜRT Academy,

Senior Consultant: Neuron Solutions,

Chief organizer: Budapest.AI."

CONTACT

# THE CHALLENGE



### Dataset

## COVID-19 Open Research Dataset Challenge (CORD-19)

An AI challenge with AI2, CZI, MSR, Georgetown, NIH & The White House

2227

**Ai2** Allen Institute For AI and 8 collaborators • updated 5 days ago (Version 3)

https://www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge

Challenge link:



This presentation:

# THERE IS UP TO DATA CASE DATA AVAILABLE

VISUALIZATION BY JHU:  https://coronavirus.jhu.edu/map.html

DATA ON GITHUB:

https://github.com/CSSEGISandData/COVID-19

BIT HIGH LEVEL

# PROBLEM FORMULATION

- There is a huge amount of literature being produced in a stream about COVID-19 and it's broader context.

- From a scientific perspective, the rapid pace of produced information can easily be overwhelming.

- We, from the AI and NLP community should lend a helping hand by mining knowledge from research papers and answer specific task questions.



DETAILS:
Challenge description

# "COMPETITION"



- Kaggle spirit, open and cooperative

- Prize is 1000USD per task, but <span style="color:red">can be offered for charity</span>

- Very few restrictions apply (eg. over 18 years of age)

- Tasks are evaluated individually, according to task specific (somewhat vague) criteria (understandably)

DETAILS:
Challenge description

PILE OF TEXT + METADATA

# THE DATASET

All in all 29000 articles + metadata

Full text, in json: **13202**

"Majority of files are in json format.

The files are grouped in 4 folders and 4 tar archives."

Approx. 2GB JSON

Pretty decent corpus with ~4.5M (?) tokens, ~385k vocabulary, but beware, <span style="color:red">multiple domains</span>!

DETAILS:
Challenge description

One notebook with some descriptives

# THE DATASET

```
paper_id :

metadata :                    dict
    title :
    authors :                 list (!)

abstract :                    list (!)

body_text :                   list (!)
    text :
    cite_spans :              list (!)
    ref_spans :               list (!)
    section :

bib_entries :                 dict (???)
    BIBREF0 :

ref_entries :                 dict (???)
    FIGREF0 :

back_matter :                 list (!)
    text :
    cite_spans :              list (!)
    ref_spans :               list (!)
    section :
```

**Observations:**

- Metadata, like author, citation (citation graph can be built)

- Text is not in one chunk, needs to be merged

- Also rich reference information is mapped "onto" the text

- One has to pay attention to dict vs. list usage

DETAILS:
Challenge description

One notebook with some descriptives

Notebook with a citation network iomplementation

VERY BROAD QUESTIONS!

# THE TASKS

## 10 "Tasks" in total:

1. What is known about transmission, incubation, and environmental stability?
2. What do we know about COVID-19 risk factors?
3. What do we know about virus genetics, origin, and evolution?
4. Sample task with sample submission - Help us understand how geography affects virality.
5. What do we know about non-pharmaceutical interventions?
6. What do we know about vaccines and therapeutics?
7. What do we know about diagnostics and surveillance?
8. What has been published about information sharing and inter-sectoral collaboration?
9. What has been published about ethical and social science considerations?
10. What has been published about medical care?

## Observations:

- Very broadly defined tasks!

  - Some hint towards very specific propositions (or entities, relations) to mine for (Like eg. 6)

  - Some are more open ended, general, like "What has been said about...", hint to "summaries" or textual parts (Like eg. 8)

- No easy mapping between taks and NLP methods, so requires thought to re-formulate task as well as to reason about the results

DETAILS:

Task definitions

# THE MAIN CHALLENGES - SPECIALIZED DOMAIN(S!)

**Observations:**

- Specialized domain(s!!!)

    - Different from general language, so out of the box external resources have to be adapted (think word vectors, or something like SciBERT?)

    - Specialized taxonomy, complex Multi-Word Expressions
        - If available, taxonomical resources have to be adapted, but maybe even unavailable
        - This is maybe the most interesting data, so good Multi-Word Expression / keyword handling is of importance

- Multiple domains, not just strict epidemIology!

## Summary statistics:

| Dataset | # Articles |
|---|---|
| CORD-19 | 29500 |
| After 2020 | 1687 |
| Uniques | 1523 |
| Covid-19 | 913 |
| Covid-19 has_full_text | 206 |

DETAILS:

Notebook source for the table right
Notebook with SciBERT embeddings

# THE MAIN CHALLENGES - SPECIFIC ANSWERS?

**Observations:**

- It can be, that we are mining for "specifics"

- Single propositions may be of strong interest

("The administration of ...(drug)... had significant effect.")

- Sometimes these may come from tabular like parts of the text

(Have to be checked!)

- Parameters and numeric values can also be of interest ("X dose of Y")



DETAILS:

One notebook with NER training in SpaCy

# THE MAIN CHALLENGES - OTHER CONSIDERATIONS

**Observations:**

- Merging of the citation graph information with the textual can be valuable

(think:weighting of a texts's importance?), but non-trivial

- The <span style="color:red">temporal order</span> of incoming information is highly relevant,

it can be, that later findings override the validity of previous propositions.

- Quantification of <span style="color:red">uncertainity</span> regarding information can be of high value

# POSSIBLE AVENUES OF "ATTACK"

**Directions it might be worth pondering:**

- Topic mining (+topic changes in time?)

- Text summarization (extractive?) methods for a relevant topic

- Search solutions or Q&A models with supporting evidence

    - Domain adapted semantic vectors

- Visualization of co-occurrence graph and frequency RELEVANT information

- Custom trained Named Entity Recognizer on entities of interest

- Multi-Word Expression detection and keyword extraction techniques

- Knowledge Graph Mining, Fact Extraction and it's relevant versions
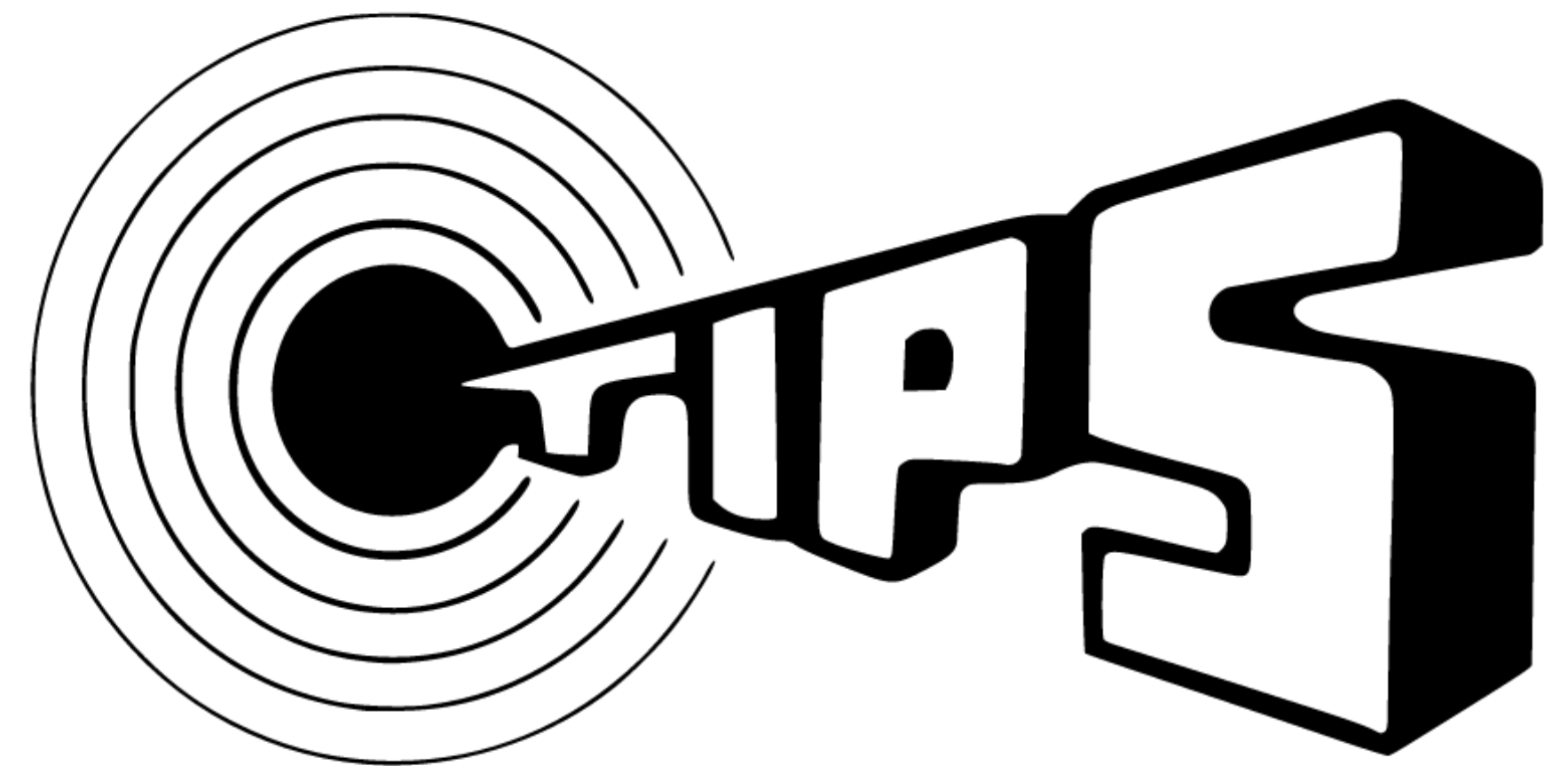
# WHY PARTICIPATE?

**- Learn (yourself and from others) & practice**

**- A chance to learn about different industries**

**(not just data science)**

**- Be part of the global community**

**- A perk to your CV or even get hired**

**- Get some cash too!**

kaggle

# TIPS FOR BEGINNERS

- **Set incremental goals**

- **Review most voted kernels**

- **Asks questions on the forums**

- **Work solo to develop core skills**

- **Team up to learn more from others**

- **Don't worry about low ranks**

Source

# NEXT STEPS:

**- Read collected papers, repos and contribute!***

**- Join our [Slack]!**

**- Brainstorm!**

**- Form teams?**

**- Save the world!**



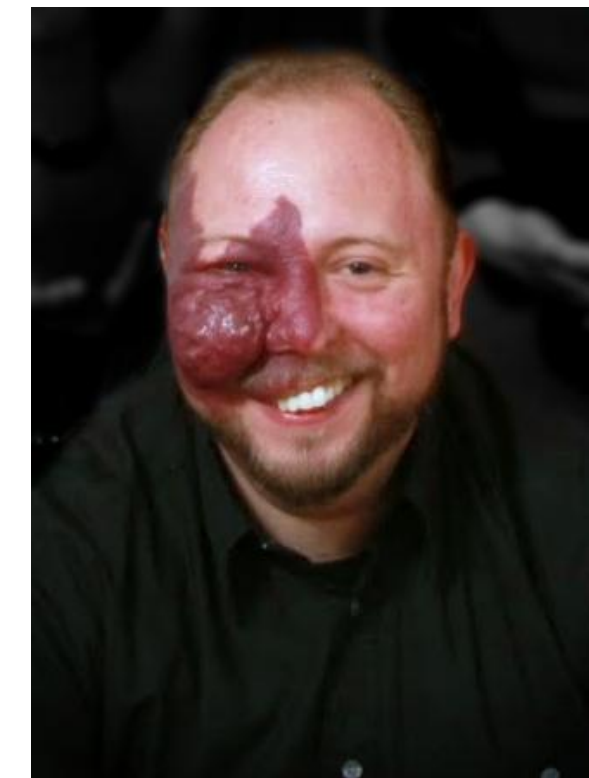*(Meta Warning! These are papers we use to process papers!)

# LET'S CONTINUE!

data-hackers-zone.slack.com

PRESENTATION

COMMUNITY

LEV

ELDAR