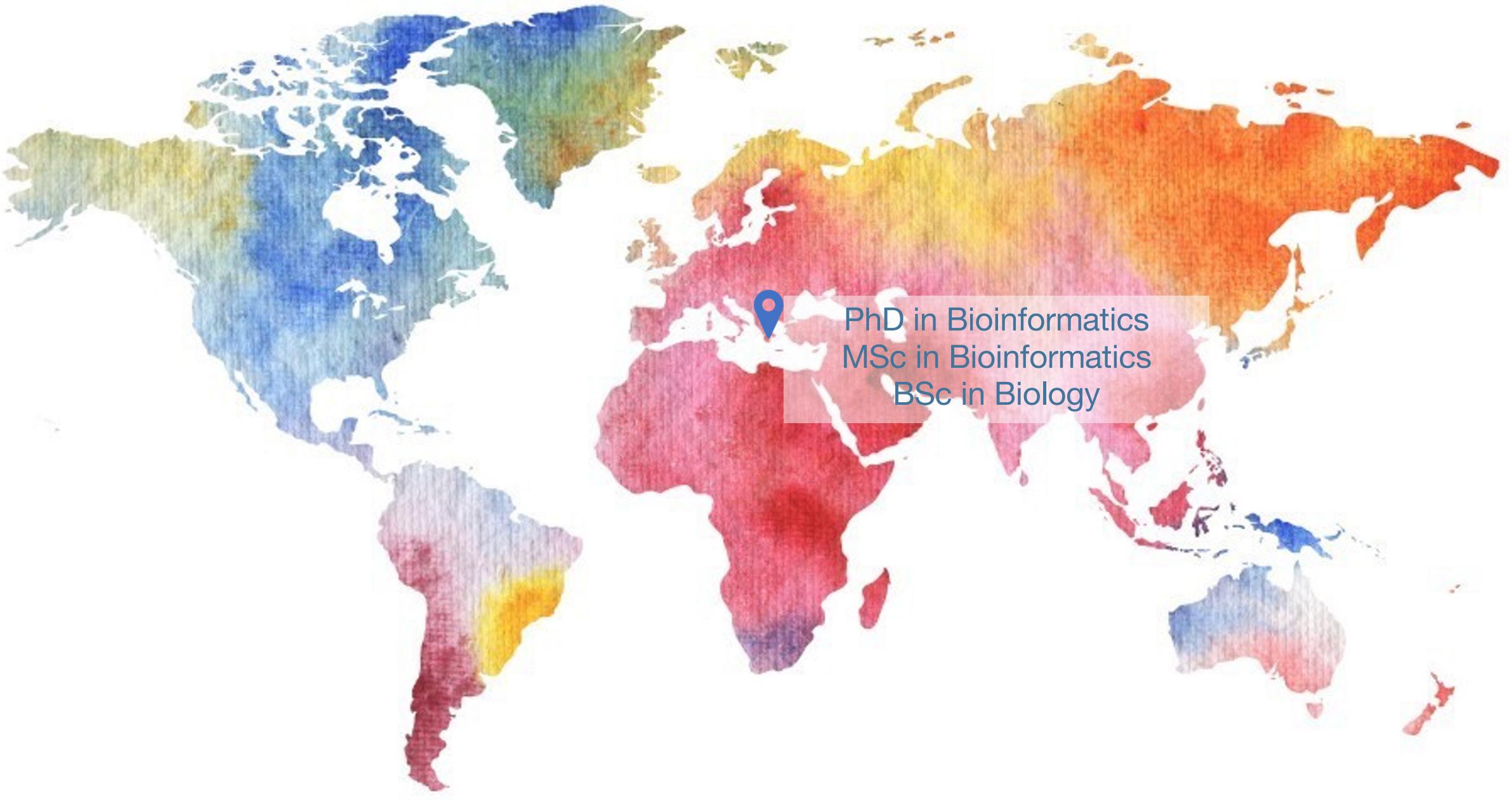


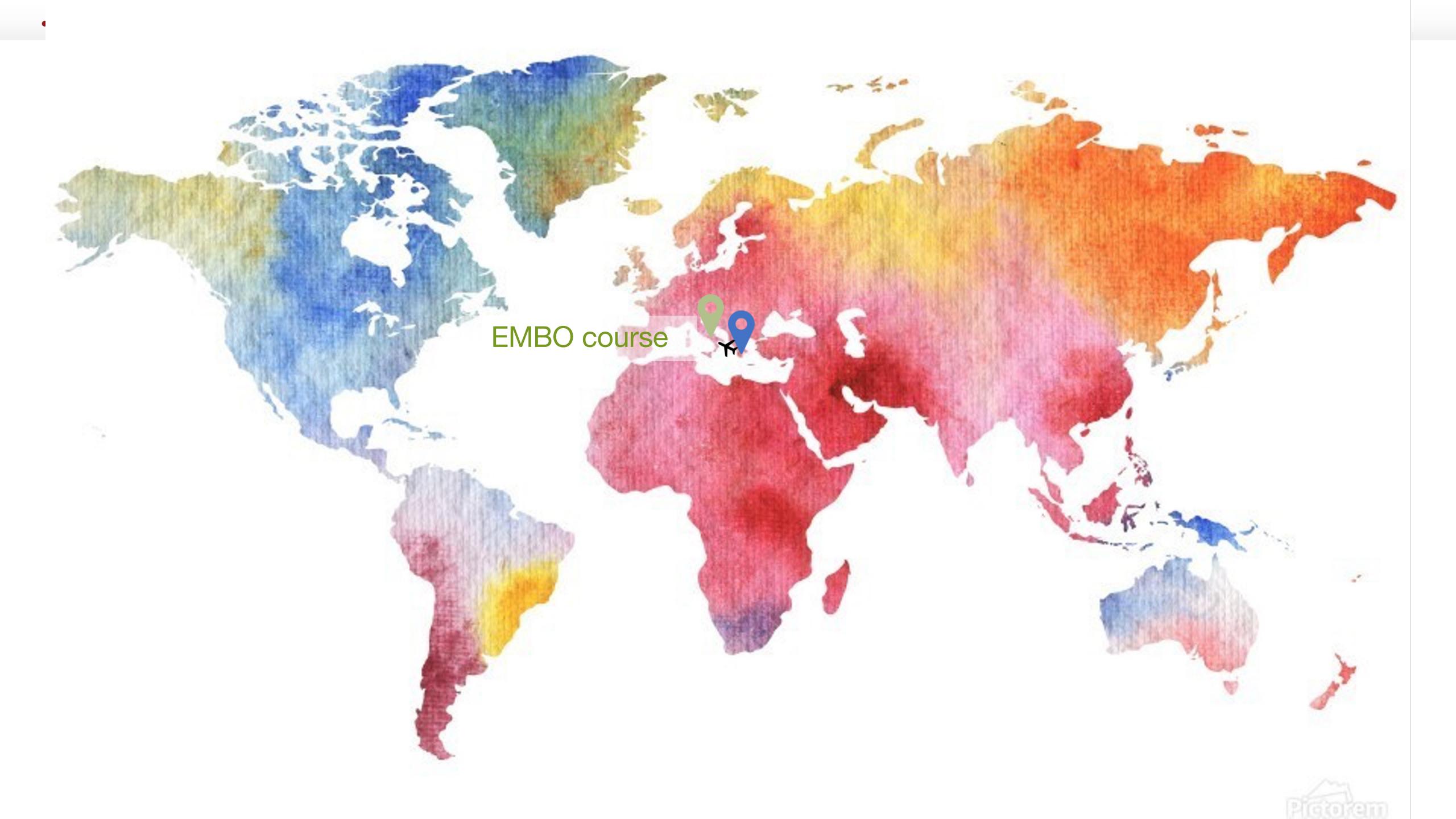
# Extracting protein-protein interactions from the literature with deep learning- based text mining

Katerina Nastou  
Cellular Network Biology Group  
Novo Nordisk Foundation Center for Protein Research  
Machine Learning course, 27-04-2023

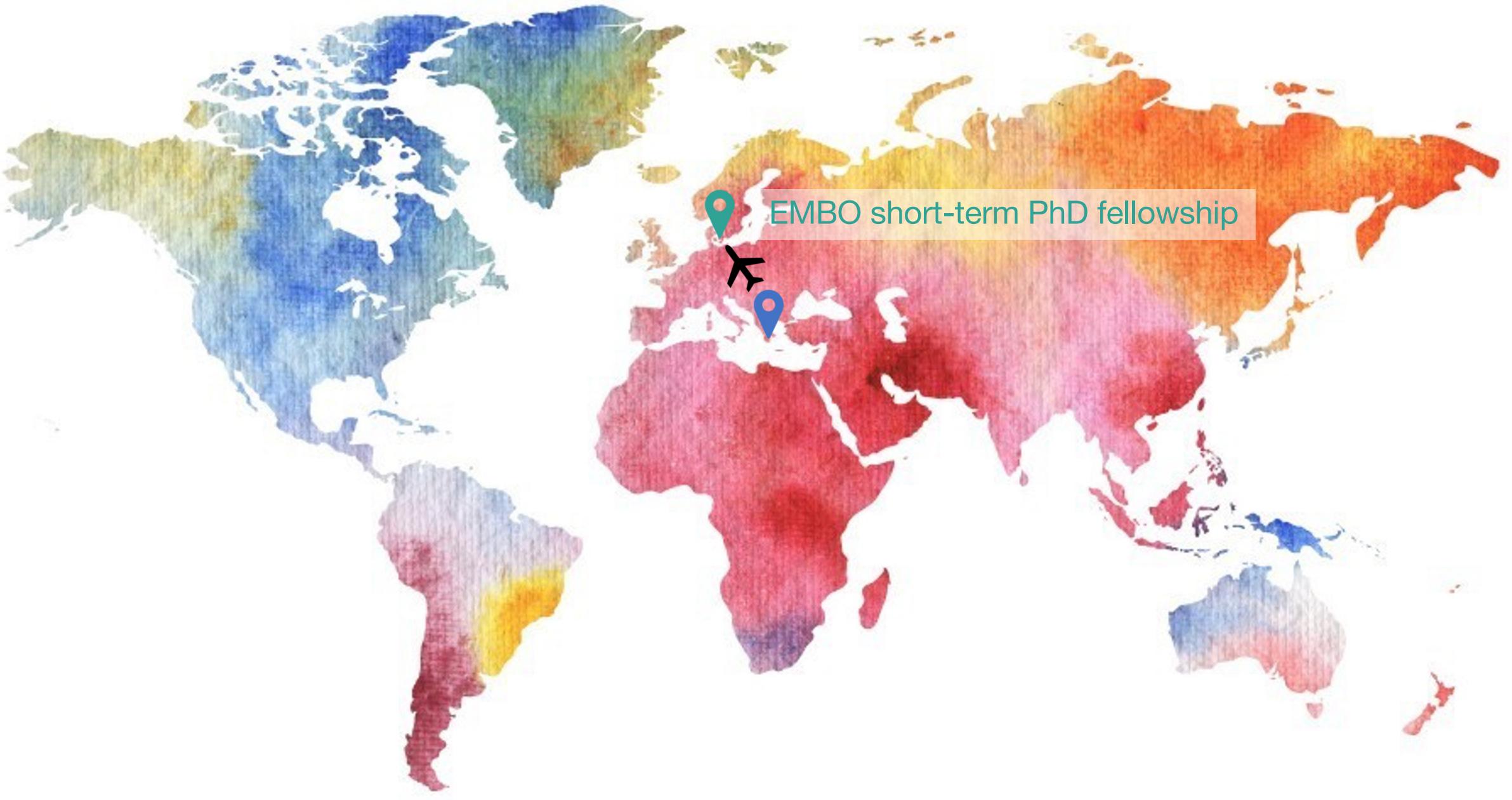
UNIVERSITY OF COPENHAGEN

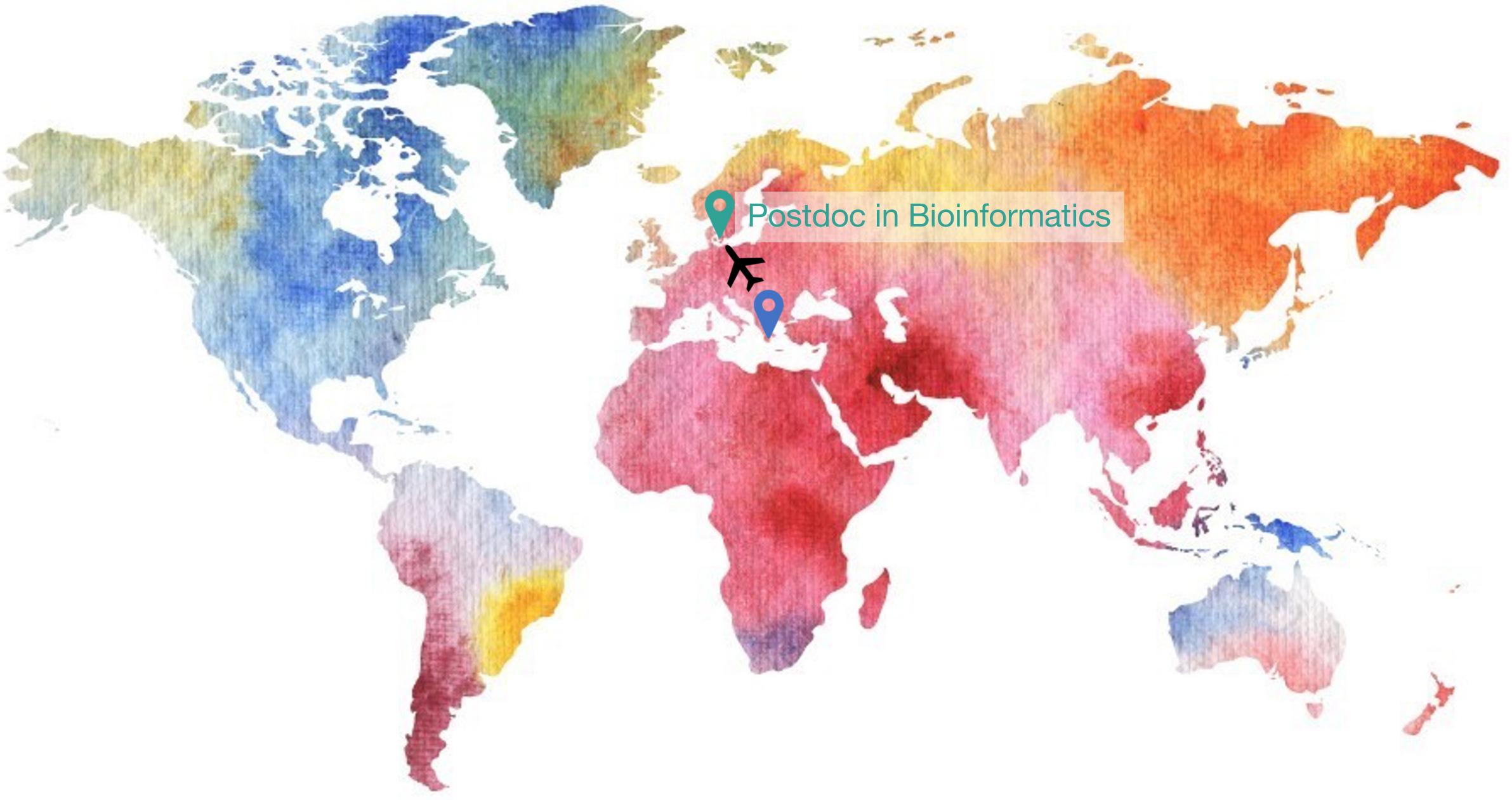






EMBO course







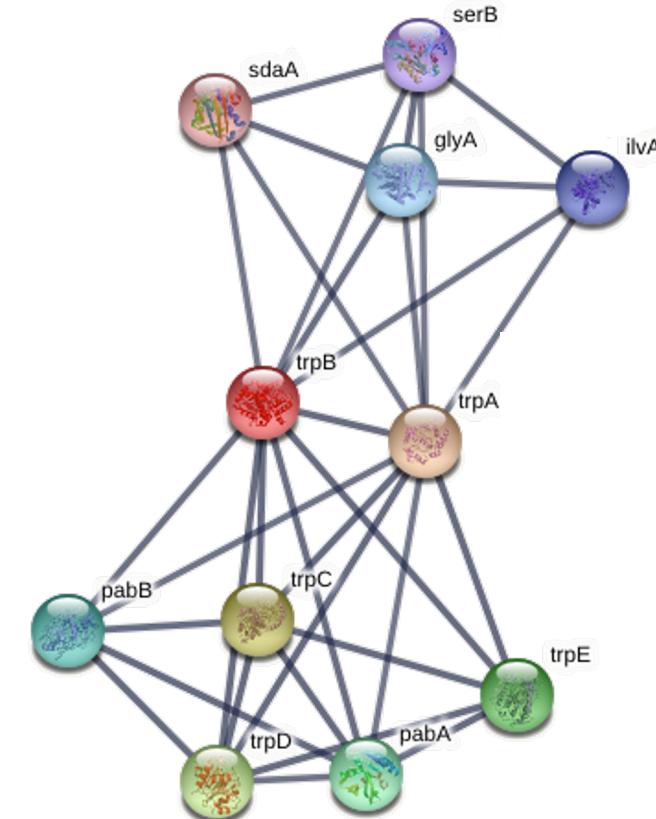


Deep learning-based text mining



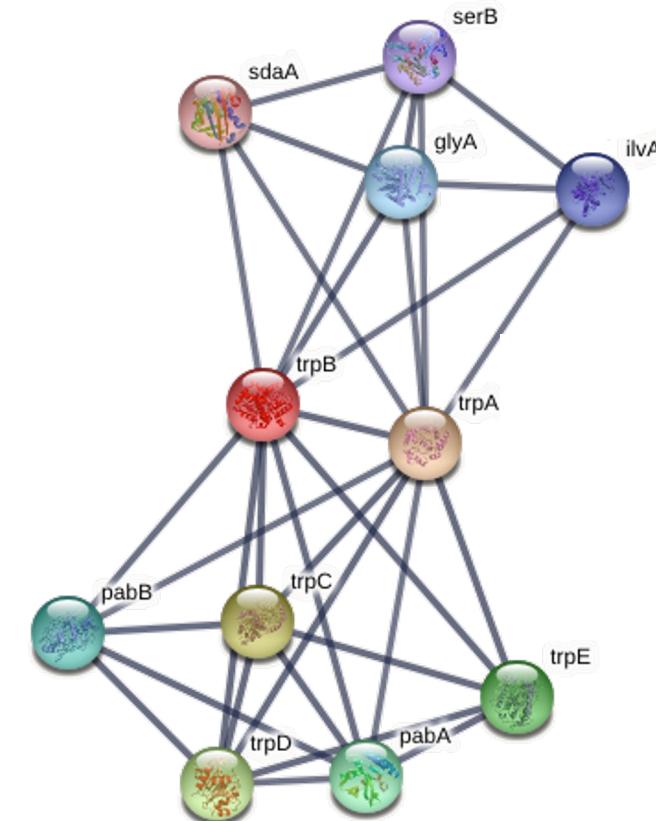
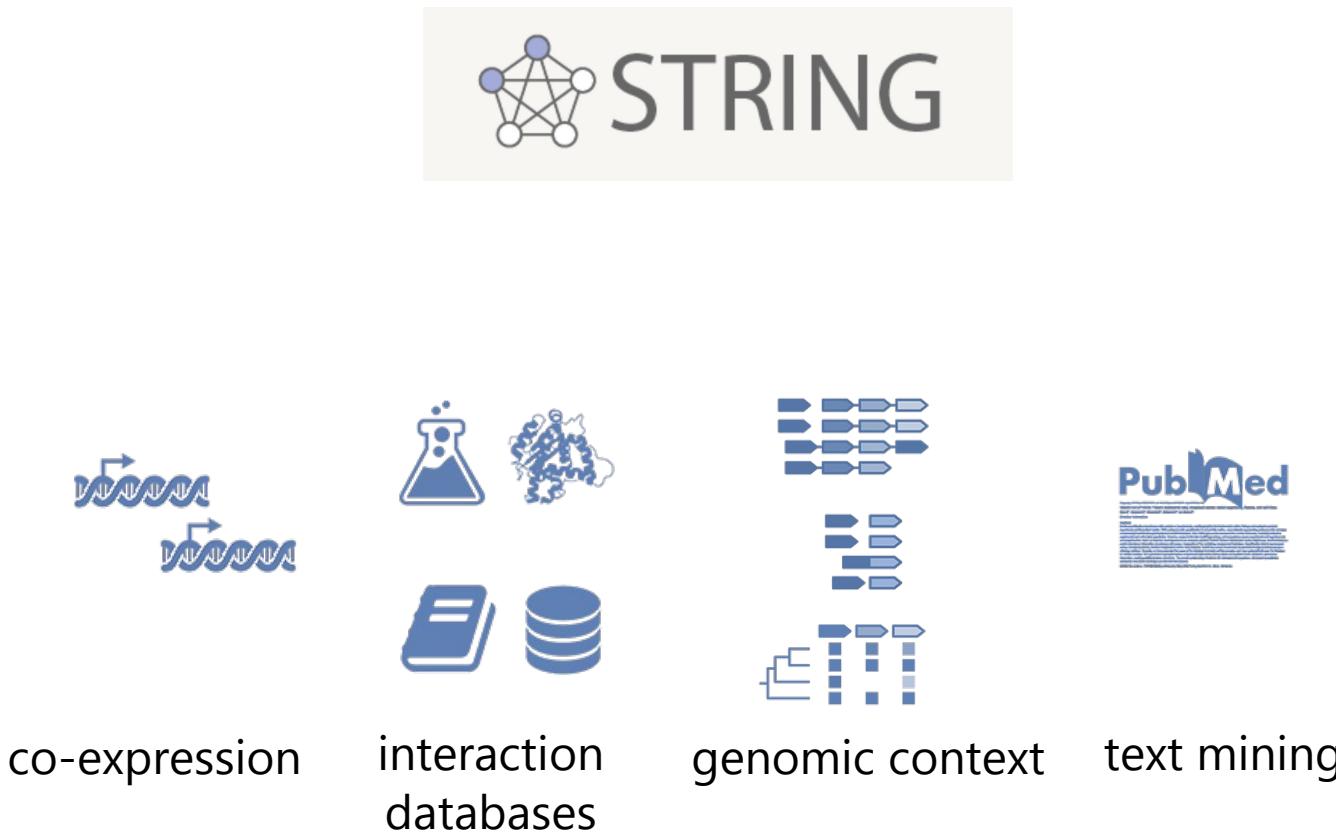
Network biology  
Protein function prediction  
Database design and development

# STRING database



Szklarczyk D, Kirsch R, Koutrouli M, Nastou K, Mehryary F, Hachilif R, Gable AL, Fang T, Doncheva NT, Pyysalo S, Bork P, Jensen LJ, von Mering C. (2022). *The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest*. Nucleic Acids Research, Jan 6;51(D1):D638-D646.

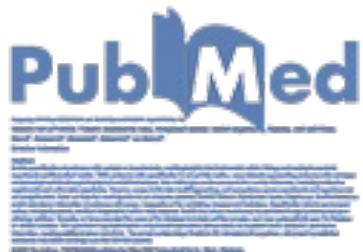
# STRING database



Szklarczyk D, Kirsch R, Koutrouli M, Nastou K, Mehryary F, Hachilif R, Gable AL, Fang T, Doncheva NT, Pyysalo S, Bork P, Jensen LJ, von Mering C. (2022). *The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest*. Nucleic Acids Research, Jan 6;51(D1):D638-D646.

~100,000 users monthly

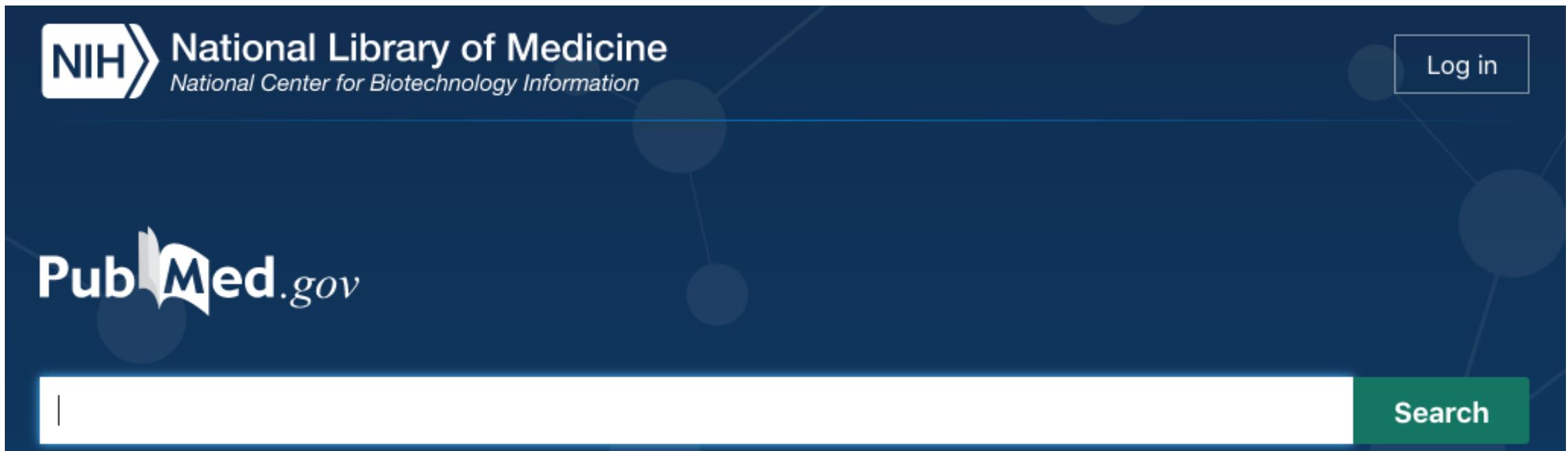
# Text Mining



# Scientific literature



# PubMed



# PubMed

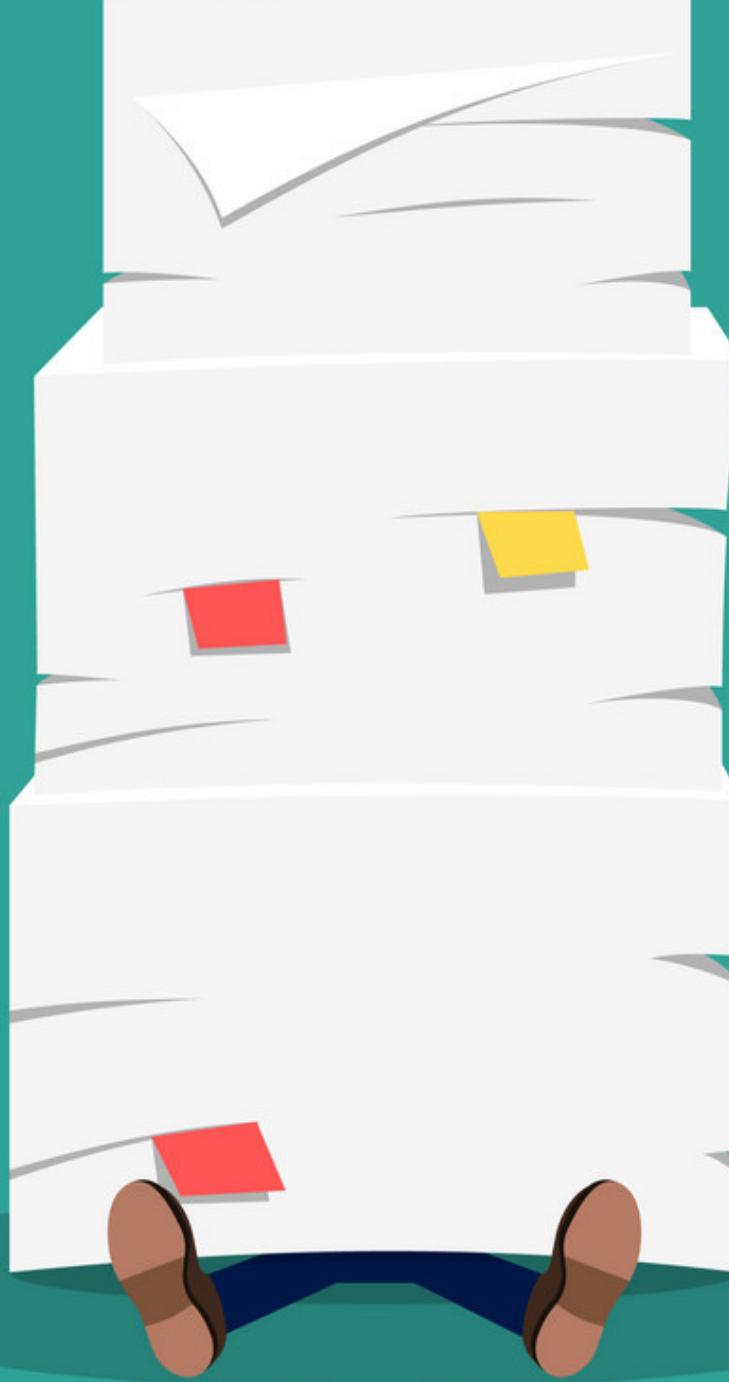
*>34 million articles  
abstracts only*

# PubMed

*~5 pages each*

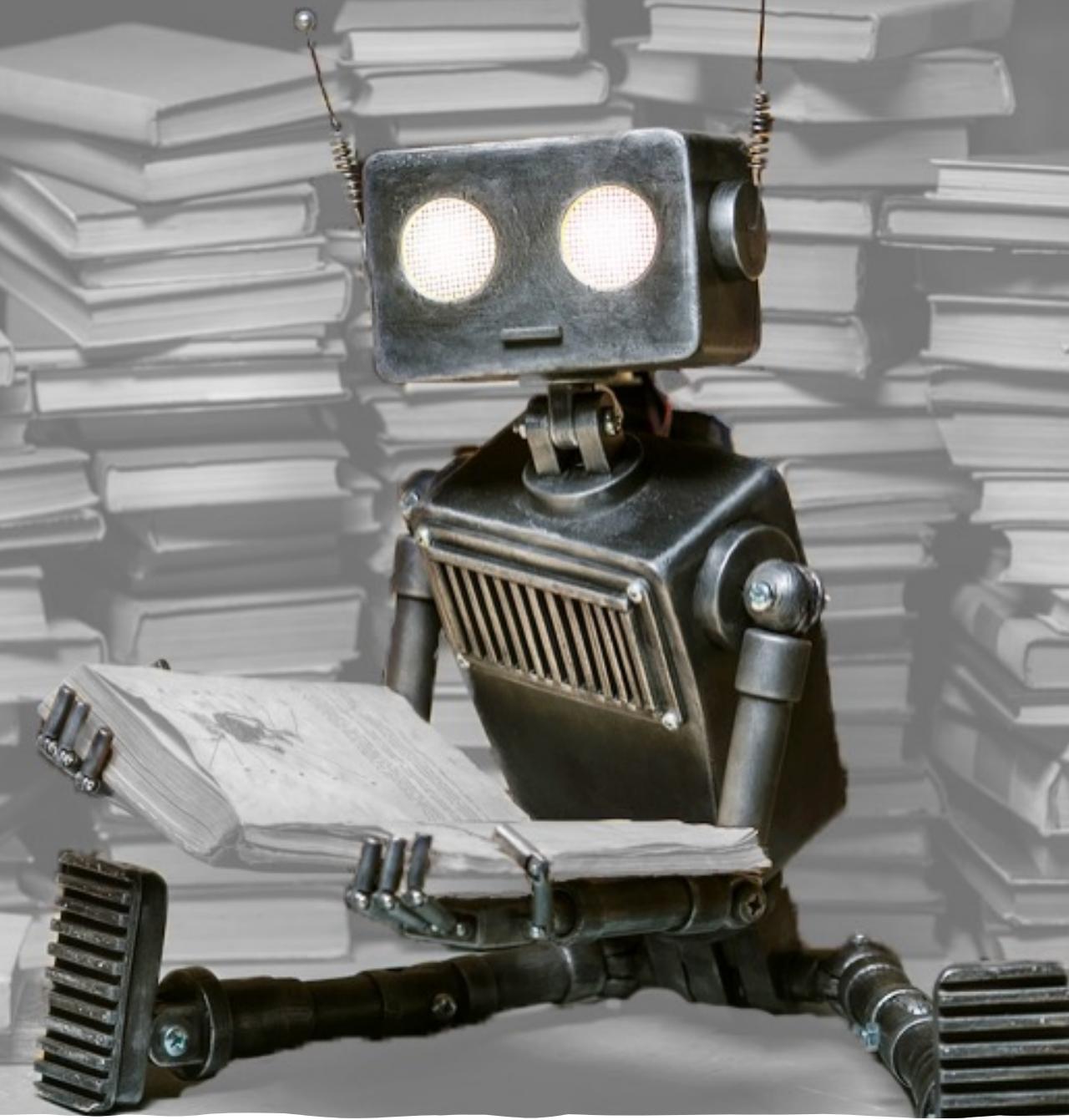
*Print in standard 80gr A4 paper*

> 10km



> 10km

Too much to read!



# STRING text-mining channel before v11.5: Statistical analysis of co-occurrence in documents + Natural Language Processing

Franceschini A., Szklarczyk D., Frankild S., Kuhn M., Simonovic M., Roth A., Lin J., Minguez P., Bork P., von Mering C., Jensen L.J. *STRING v9.1: protein-protein interaction networks, with increased coverage and integration.* Nucleic Acids Res. 2013

# STRING text-mining channel: Statistical analysis of co-occurrence in documents + Natural Language Processing

## Functional associations

Franceschini A., Szklarczyk D., Frankild S., Kuhn M., Simonovic M., Roth A., Lin J., Minguez P., Bork P., von Mering C., Jensen L.J. *STRING v9.1: protein-protein interaction networks, with increased coverage and integration.* Nucleic Acids Res. 2013

# STRING text-mining channel before v11.5: Statistical analysis of co-occurrence in documents + Natural Language Processing

physical and regulatory interactions

# STRING text-mining channel before v11.5 :

## Statistical analysis of co-occurrence in documents + Natural Language Processing



physical and regulatory interactions

# STRING text-mining channel before v11.5 :

## Statistical analysis of co-occurrence in documents + Natural Language Processing

physical and regulatory interactions



# STRING text-mining channel before v11.5 :

Statistical analysis of co-occurrence in documents +  
Natural Language Processing

physical and regulatory interactions



# STRING text-mining channel before v11.5 :

## Statistical analysis of co-occurrence in documents + Natural Language Processing

physical and regulatory interactions



# 🔥 Deep Learning-based Language Representation Models 🔥

# BERT



# What is BERT?

# BERT: Bidirectional Encoder Representations from Transformers

# BERT: Bidirectional Encoder Representations from Transformers

Language modelling system,  
pre-trained with unlabeled  
data, then fine-tuned

# BERT: Bidirectional Encoder Representations from Transformers

Language modelling system,  
pre-trained with unlabeled  
data, then fine-tuned

BERT is not made to tackle a specific task or answer a specific problem but it gives the best, the most efficient and the most flexible representation for words and sequences possible

# BERT: Bidirectional Encoder Representations from Transformers

Language modelling system,  
pre-trained with unlabeled  
data, then **fine-tuning**

Add a small layer to the model made for our specific task and train the whole thing very lightly  
(really fast, because the core of the model has been already trained)

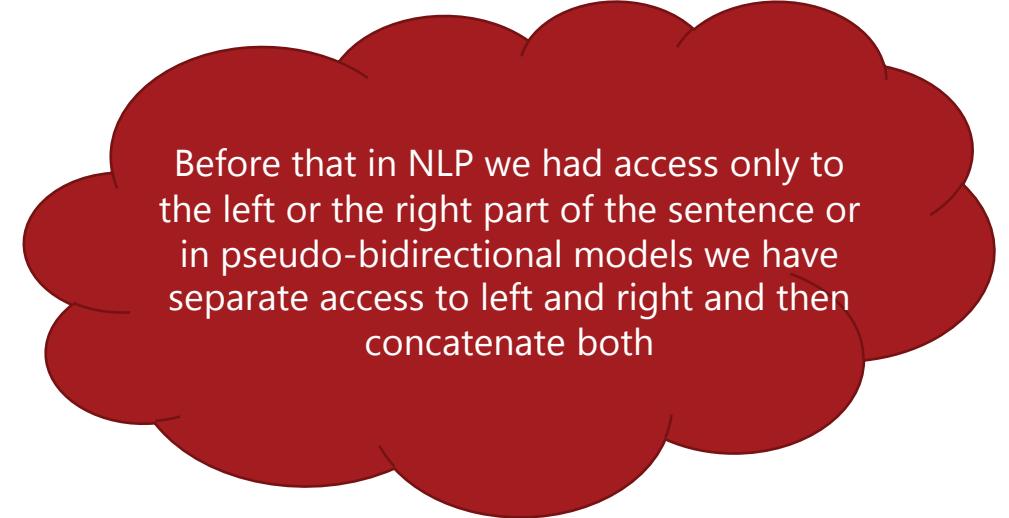
# BERT: Bidirectional Encoder Representations from Transformers

Defines the powerful neural network architecture of BERT, that is based on a self-attention mechanism

# BERT: Bidirectional Encoder Representations from Transformers

Has full access to left  
and right **context**  
when dealing with a  
word (masked)

Defines the training  
process

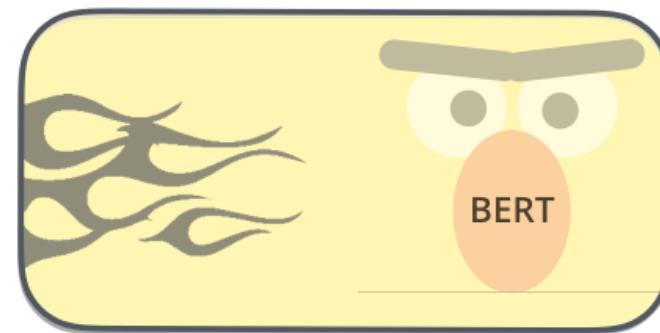


Before that in NLP we had access only to the left or the right part of the sentence or in pseudo-bidirectional models we have separate access to left and right and then concatenate both

## 1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

### Semi-supervised Learning Step



Model:



WIKIPEDIA  
Die freie Enzyklopädie

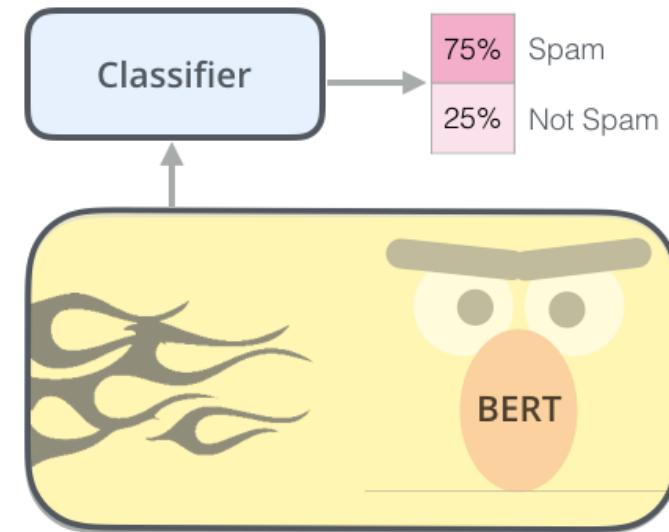
Dataset:

Predict the masked word  
(language modeling)

Objective:

## 2 - Supervised training on a specific task with a labeled dataset.

### Supervised Learning Step



Model:  
(pre-trained  
in step #1)

Dataset:

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

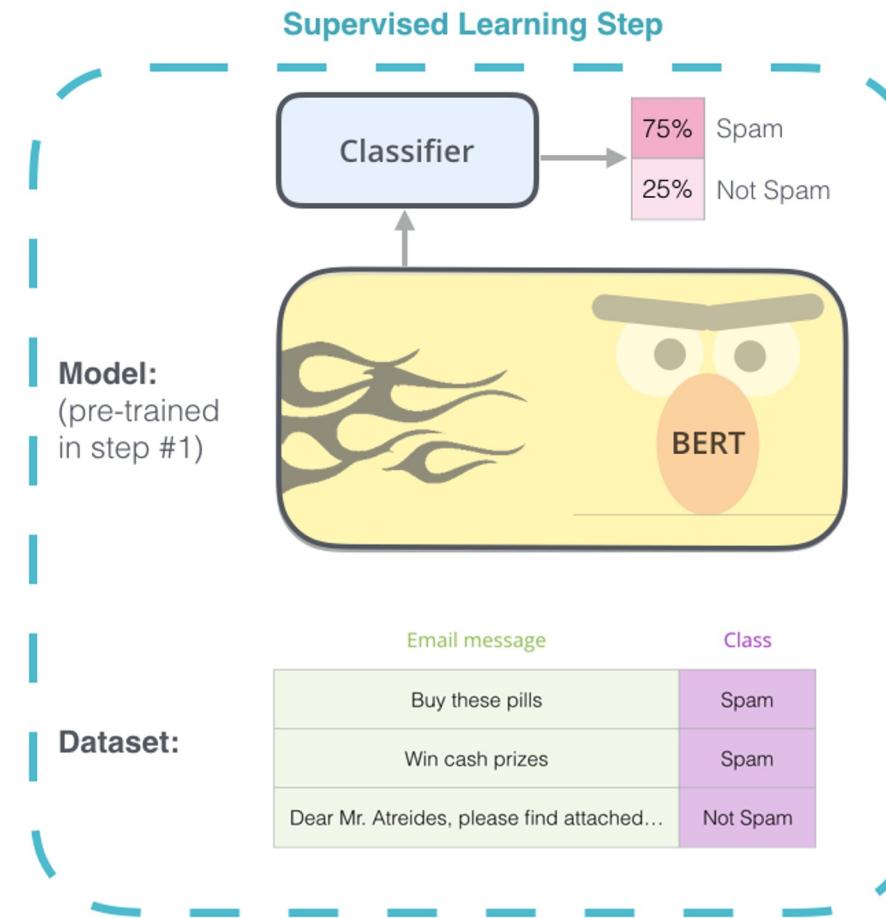
# BERT: Bidirectional Encoder Representations from Transformers

During **supervised training** words of interest are **masked** and the model learns the label for the masked word based on the **context** around them

Example	Label
When [MASK] becomes mutated, it loses its function, resulting in abnormal cell proliferation and tumor progression.	gene BERT
Patients suffering from [MASK] have a low white blood cell count (WBC).	disease WIKIPEDIA Die freie Enzyklopädie

Predict the masked word (language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.

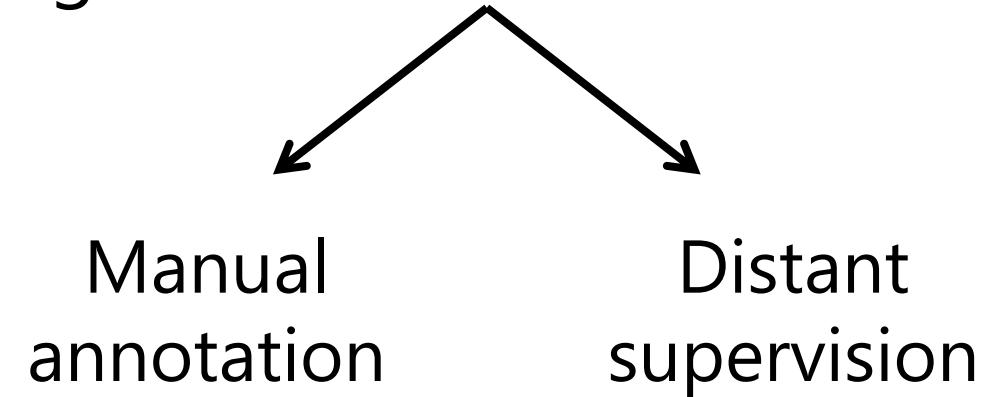


# Domain-specific deep-learning models

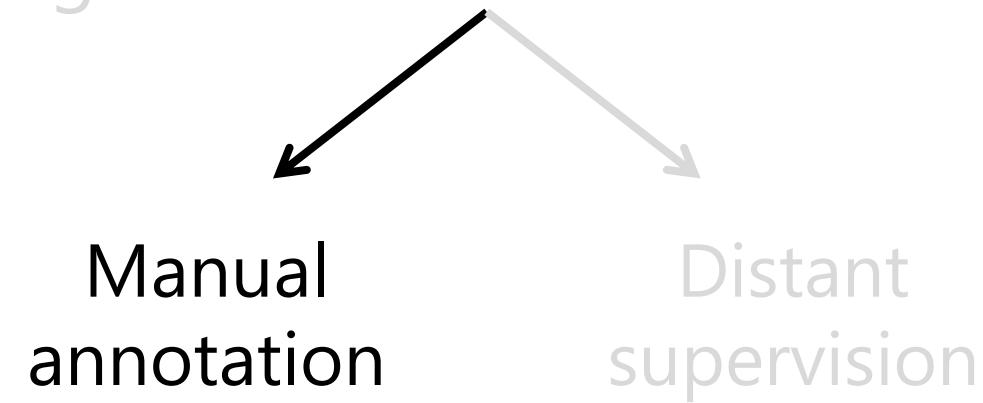
Domain-specific deep-learning models  
BioBERT, SciBERT, BlueBERT (NCBI-BERT),  
BioMegatron, RoBERTa-bio and more...

Major challenge in using deep learning for Relation Extraction  
Obtaining enough amounts of labelled text for training

Major challenge in using deep learning for Relation Extraction  
Obtaining enough amounts of labelled text for training



Major challenge in using deep learning for Relation Extraction  
Obtaining enough amounts of labelled text for training



## Manual Annotation

- +good quality
- +more detailed annotations
- manual

# Development of guidelines for manual annotation

<https://katnastou.github.io/stringdb-typed-relation-annotation-docs>



And there is more...



<https://brat.nlplab.org/>

/string-relation-corpus/physical-interaction-dbs-abstracts-01/10021333 brat

1 Apontic binds the translational repressor Bruno and is implicated in regulation of oskar mRNA translation.

2 The product of the oskar gene directs posterior patterning in the Drosophila oocyte, where it must be deployed specifically at the

3 Proper expression relies on the coordinated localization and translational control of the oskar mRNA.

4 Translational repression prior to localization of the transcript is mediated, in part, by the Bruno protein, which binds to discrete sit

5 To begin to understand how Bruno acts in translational repression, we performed a yeast two-hybrid screen to identify Bruno-inte

6 One interactor, described here, is the product of the apontic gene.

7 Coimmunoprecipitation experiments lend biochemical support to the idea that Bruno and Apontic proteins physically interact in

8 Genetic experiments using mutants defective in apontic and bruno reveal a functional interaction between these genes.

9 Given this interaction, Apontic is likely to act together with Bruno in translational repression of oskar mRNA.

10 Interestingly, Apontic, like Bruno, is an RNA-binding protein and specifically binds certain regions of the oskar mRNA 3' untranslated region.

**Edit Annotation**

**From**  
Gene or gene product ("Apontic") [Link](#)

**To**  
Gene or gene product ("oskar")

**Type**

- Regulation
  - Positive regulation
  - Negative regulation
- Complex formation
- Regulation of gene expression
  - Regulation of translation
  - Regulation of transcription
- Regulation of proteolysis
- Catalysis of protein modification
  - Catalysis of phosphorylation

**Notes**

[Reverse](#) [Delete](#) [Reselect](#) [OK](#) [Cancel](#)

Identify relationships based on context

What relationships?

v11.5: Identify **physical interactions** based on context

... interaction of  
A and B ...

... A-B complex ...

## Identify physical interactions based on context

... A interacts  
with B ...

... A binds to B ...

... A/B  
heterodimer ...

# Where to get textual information from?

Physical interactions (BioGRID, IntAct, MINT)

Previous BioNLP shared tasks annotation datasets

Oughtred R, Stark C, Breitkreutz BJ, Rust J, Boucher L, et al. *The BioGRID interaction database: 2019 update*. Nucleic Acids Research, 2019

Orchard S, Ammari M, Aranda B, Breuza L, Brigandt I, et al. *The MIntAct project--IntAct as a common curation platform for 11 molecular interaction databases*. Nucleic Acids Research, 2014

Pyysalo S, Ohta T, Rak R, Rowley A, Chun HW, Jung SJ, et al. *Overview of the cancer genetics and pathway curation tasks of BioNLP shared task 2013*. BMC Bioinformatics 2015

Kim JD, Pyysalo S, *BioNLP Shared Task*, Encyclopedia of Systems Biology, 2013

## Corpus (version 1)

Physical interactions (BioGRID, IntAct, MINT)

400 documents (300 abstracts + 100 full-text paragraphs)

Previous BioNLP shared tasks annotation datasets

137 documents enriched in complex formation events

## Training dataset (version 1)

Physical interactions (BioGRID, IntAct, MINT)

400 documents (300 abstracts + 100 full-text paragraphs)

Previous BioNLP shared tasks annotation datasets  
137 documents enriched in complex formation events

Total no of relationships:

**Positive** (complex formation): 1607

**Negative** (Not a complex): 10392

## Training dataset (version 1)

Total no of relationships:

Positive (complex formation): 1607

Negative (Not a complex): 10392



Shuffle data &  
reduce redundancy

**Positive** (complex formation): 1298

**Negative** (Not a complex): 6385

## Training dataset (version 1)

Total no of relationships:

Positive (complex formation): 1607

Negative (Not a complex): 10392



Shuffle data &  
reduce redundancy

Labelled dataset  
to fine tune  
BioBERT model

**Positive** (complex formation): 1298  
**Negative** (Not a complex): 6385

## Training dataset (version 1)

Total no of relationships:

Positive (complex formation): 1607

Negative (Not a complex): 10392



Shuffle data &  
reduce redundancy

Labelled dataset  
to fine tune  
BioBERT model

**Positive** (complex formation): 1298  
**Negative** (Not a complex): 6385

Even less  
depending  
on the max  
seq len

# Test various hyperparameters (grid search)

Models: **BioBERT\_base**, **BioBERT\_large**

Batch size: 32, 64, 256, 512, 1024, 2048

Learning rate: 5E-6, 1E-5, 2E-5

Epochs: 5, 10, 20

Maximum sequence length: 96 (large), 256 (base)

# Test various hyperparameters

## Best combination

Batch size: 32, 64, 256, 512, 1024, 2048

Learning rate: 5E-6, 1E-5, 2E-5

Epochs: 5, 10, 20

Maximum sequence length: 96 (large), 256 (base)

# Test various hyperparameters

## Best combination

Batch size: 32, 64, 256, 512, 1024, 2048

Learning rate: 5E-6, 1E-5, 2E-5

Epochs: 5, 10, 20

Maximum sequence length: 96 (large), 256 (base)

F1-score: 83.1%



## Training dataset (version 1)

Total no of relationships:

Positive (complex formation): 1607

Negative (Not a complex): 10392

Labelled dataset  
to fine tune  
BioBERT model



Shuffle data &  
reduce redundancy

Positive (complex formation): 1298  
Negative (Not a complex): 6385

The culprit?

## Use Best model

Run Prediction on all pairs of proteins in the literature  
that are in the same sentence (large scale run)

Benchmarking and Score calibration for physical interactions  
against **Complex Portal**

Physical interactions mode for version 11.5 of STRING

Check our paper for more details



Szklarczyk, D., Gable, A. L., Nastou, K. C., Lyon, D., Kirsch, R., Pyysalo, S., Doncheva, N. T., Legeay, M., Fang, T., Bork, P., Jensen, L. J., and von Mering, C. (2021). *The STRING database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets*. Nucleic Acids Res, 49(D1), D605-612

v12: Identify **physical interactions** based on context (again)

## Corpus (version 2)

Physical interactions (BioGRID, IntAct, MINT)

~~400 documents (300 abstracts + 100 full-text paragraphs)~~  
800                  400                  400

Previous BioNLP shared tasks annotation datasets

137 documents enriched in complex formation events

Regulatory Interactions (Reactome)

300 documents (abstracts)

Genetic Interactions (BioGRID)

50 documents (abstracts)

Total number of documents in Corpus v2: 1,287

# Introduction of new Named Entity types

Gene or Gene Product (Protein → GGP)

Protein Family or Group

Protein-containing Complex

Chemical

Out-of-scope

# New training/dev/test split

## New train/dev/test split

Document-based and not sentence-based

60% train / 20% dev / 20% test doc split  
instead of 70% train / 10% dev / 20% test



avoid data leakage to dev set

60/20/20 split also applies to the complex formation relationships (not only the documents)

## Total number of positive (complex formation) relationships in corpus v2

Training set: 2,131

Dev set: 650

Test set: 644



Total number of complex formation relationships  
in training set?

Total number of complex formation relationships?  
Depends on experimental setup

Exp	Training pairs (Positive/Negatives)	Masked entities	Not masked entities	Comment	Positive pair count in train	Negative pair count in train
1-A	(X, Y) X, Y ∈ {Prot, Prot_BL, Complex, Chemical, Family}		Prot, Prot_BL, Complex, Chemical, Family	maximum training data No masking	2,489	201,673
1-B	(X, Y) X, Y ∈ {Prot, Prot_BL, Complex, Chemical, Family}	Prot, Prot_BL, Complex, Chemical, Family		maximum training data maximum masking	2,489	201,673
2	(Prot, Prot)	Prot	Prot_BL, Complex, Chemical, Family	minimum training data minimum masking most similar to clean dev	1,961	113,527
3	(X, Y) X, Y ∈ {Prot, Prot_BL, Complex}	Prot, Prot_BL, Complex, Chemical, Family		more training data compared to Experiment #2 maximum masking	2,117	138,961
4	(X, Y) X, Y ∈ {Prot, Prot_BL, Complex}	Prot, Prot_BL, Complex,	Chemical, Family	more training data compared to Experiment #2 less masking compared to Experiment #3	2,117	138,961

## Test A LOT of hyperparameters

Models: **BERT\_X**

Batch size: 16

Learning rate: 2e-5 3e-5 5e-5 5e-6 5e-7

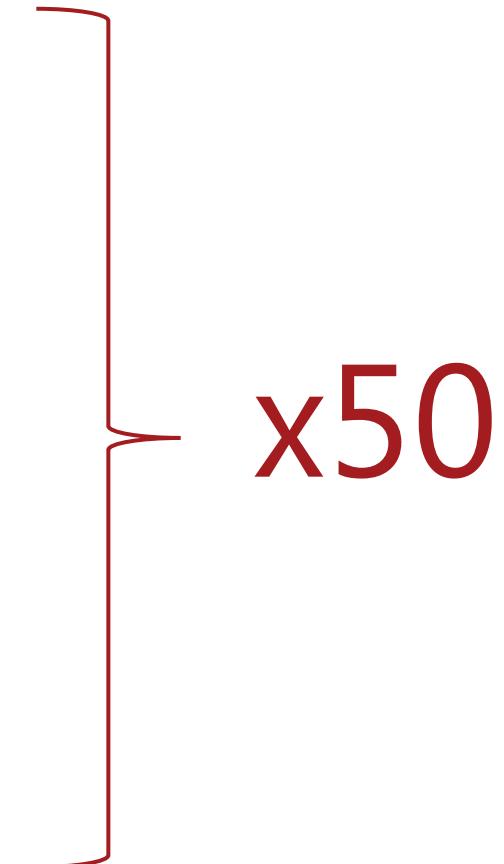
Epochs: 2, 3, 5, 8, 10, 12

Maximum sequence length: 64, 128, 256, 512

Random seed index: 0, 1, 2, 3

Optimizer: adam, adam\_warmup

Strategy: MARK, MASK



## Test various hyperparameters (final grid search)

Models: BioBERT\_base, **RoBERTa-large-PM-M3-Voc-hf**

Experimental setup: 1-A, **1-B (MASK)**, 2, 3 , 4

Batch size: 5

Learning rate: **3E-6**, 4E-6, 5E-6

Epochs: 7, 8, 9, 10, **11**, 12

Maximum sequence length: **128**, 144, 160, 176, 192

**Mean F1-score: 84.25%**  
(std = 0.82%)

Not comparable with  
previous F-score

## Best model

Run Prediction on all pairs of proteins in the literature  
that are in the **same paragraph** (large scale run)

Benchmarking and Score calibration for physical interactions  
against **Complex Portal** and for functional against KEGG

Physical interactions for version 12 of STRING

Check out the latest paper



Szklarczyk D, Kirsch R, Koutrouli M, Nastou K, Mehryary F, Hachilif R, Gable AL, Fang T, Doncheva NT, Pyysalo S, Bork P, Jensen LJ, von Mering C. (2022). *The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest.* Nucleic Acids Research, Online ahead of print

v12.5: Identify **physical interactions and typed directed relationships** based on context

# What relationships?

# Relationship types

Regulation

- Positive regulation
- Negative regulation

Complex formation

Regulation of gene expression

- Regulation of transcription
- Regulation of translation

Regulation of degradation

Catalysis of Post-translational modification

- Catalysis of small protein conjugation/removal
- Catalysis of small protein conjugation

    Catalysis of Ubiquitination

    Catalysis of SUMOylation

    Catalysis of Neddylation

    Other catalysis of small protein conjugation

- Catalysis of small protein removal

    Catalysis of Deubiquitination

    Catalysis of DeSUMOylation

    Catalysis of Deneddylation

    Other catalysis of small protein removal

Catalysis of Post-translational modification

Catalysis of phosphoryl group conjugation/removal

Catalysis of Phosphorylation

Catalysis of Dephosphorylation

Catalysis of small molecule conjugation/removal (excluding phosphoryl group)

Catalysis of small molecule conjugation

Catalysis of Methylation

Catalysis of Acylation

Catalysis of Acetylation

Catalysis of Palmitoylation

Catalysis of Myristoylation

Catalysis of lipidation

Catalysis of prenylation

Catalysis of farnesylation

Catalysis of geranylgeranylation

Catalysis of Glycosylation

Catalysis of ADP-ribosylation

Other catalysis of small molecule conjugation

Catalysis of small molecule removal

Catalysis of Demethylation

Catalysis of Deacylation

Catalysis of Deacetylation

Catalysis of Depalmitoylation

Catalysis of Deglycosylation

Other catalysis of small molecule removal

## Corpus (version 3)

Physical interactions (BioGRID, IntAct, MINT)

800 documents (400 abstracts + 400 full-text paragraphs)

Previous BioNLP shared tasks annotation datasets

137 documents enriched in complex formation events

Regulatory Interactions (Reactome)

300 documents (abstracts)

Genetic Interactions (BioGRID)

50 documents (abstracts)

Exhaustive PTM corpus

234 documents enriched in post-translational modification events

Total number of documents in Corpus v3: 1,521

## Total number of “positive” relationships in corpus v3

Training set: 5,508 (2,362 Complex formation)

Dev set: 1,778 (741 Complex formation)

Test set: 1,735 (720 Complex formation)

# Experiments?

Initial experiments: not on all relationship types

Initial experiments: not on all relationship types  
Only on those with more than 40 relationships in training set

# Setup: from binary to multi-label

## Initial experiments with small grid search

Models: **RoBERTa-large-PM-M3-Voc-hf**

Experimental setup: 1-A (MARK)

Batch size: 5

Learning rate: 2e-6, **3E-6**, 5E-6

Epochs: 8, 9, 10, **12**

Maximum sequence length: **128**

**Best mean F1-score:**  
**71.88% (std = 0.48%)**

# Future directions

Corpus v4: Expand corpus to 2,000 documents  
*(i.e. add 479 more documents)*

Corpus v4: Add more documents for PTM relationships

Corpus v5: Expand corpus to 3,000 documents

## Corpus v5: Add more full-text excerpts

Run hyperparameter grid searches with final corpus

# Different experimental setups

## Different experimental setups

Main experimental question: how to treat relationships without enough examples to train/test?

Run predictions on all protein pairs

But before adding to STRING v12.5

New benchmark set for regulatory typed relationships

In the meantime: Create network of directed typed relationships from deep learning-based text mining

Use proteomics datasets to check whether one can infer upstream regulators based on downstream events



Strange Planet  
by Nathan Pyle

# Thank you for your attention!

## JensenLab

### Lars Juhl Jensen

Stefano Roncelli

Mikaela Koutrouli

Esmaeil Nourani

Oana Palasca

### Rebecca Kirsch

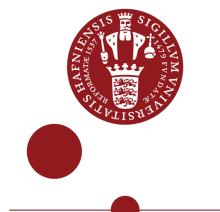
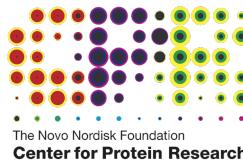
Katerina Nastou

Maud Kerstholt

Dewei Hu

Nadezhda Tsankova Doncheva

Foteini Aktypi



## TurkuNLP lab

### Sampo Pyysalo

### Farrokh Mehryary



novo  
nordisk  
fonden



Contact: [katerina.nastou@cpr.ku.dk](mailto:katerina.nastou@cpr.ku.dk)