

Exploring the versatility of LLMs for Relation extraction in underexplored biomedical areas

Maxime Delmas - 05/03/2024

Background

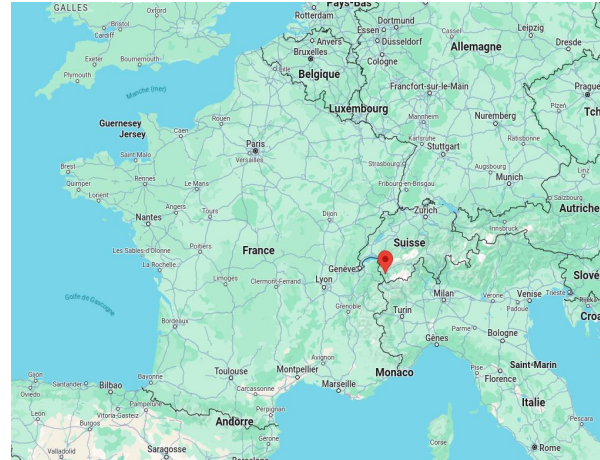


PhD (2019 - 2022)



biomedicine
Metabolomics
Bayesian statistics
Network analysis
Natural Language Processing
Knowledge Graph
Web Sémantique
Ontologies
Data Integration
Deep learning
Explainability

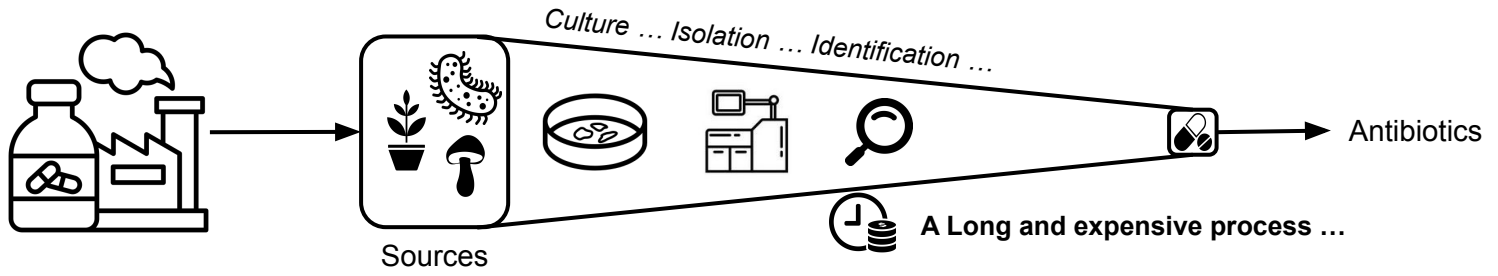
Post doctoral position (2022 - now)



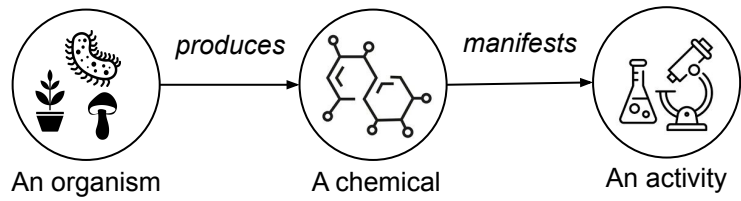
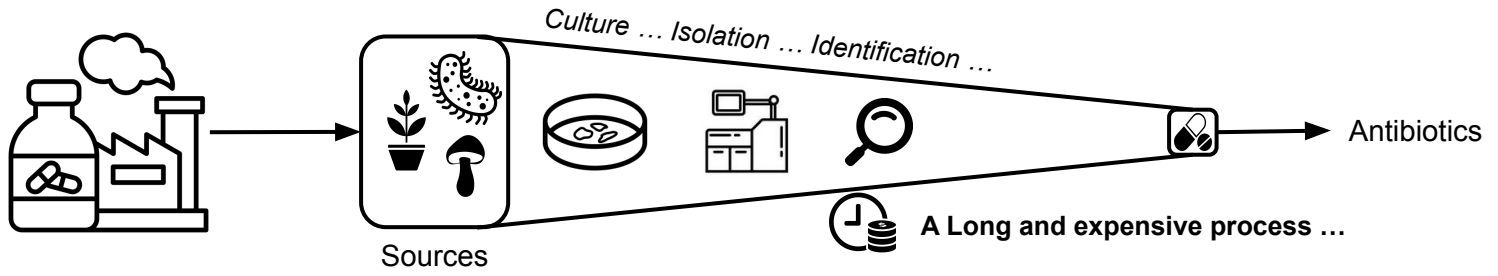
INFLAMALPS

Network analysis
Reasoning
Deep learning
Explainability
Natural Language Processing
Knowledge Graph
biomedicine
Web Sémantique
Data Integration

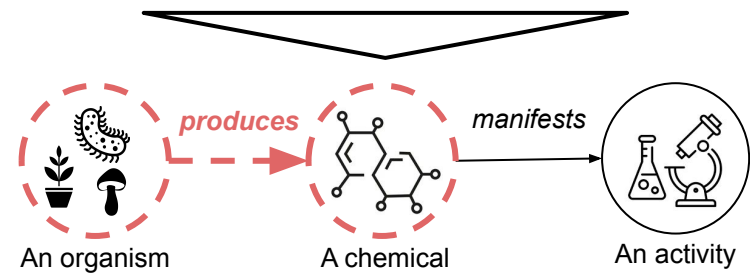
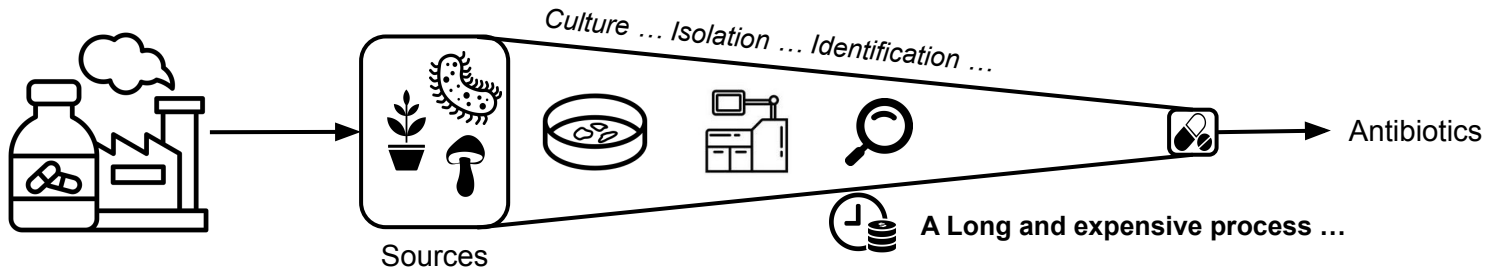
Study context: avoiding rediscoveries



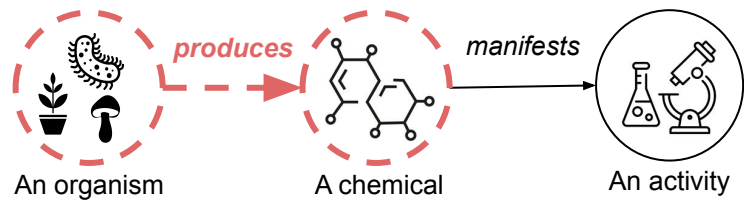
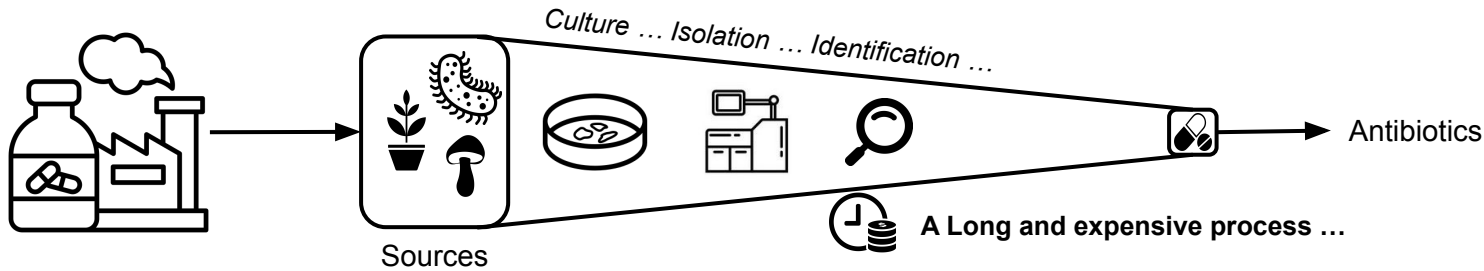
Study context: avoiding rediscoveries



Study context: avoiding rediscoveries



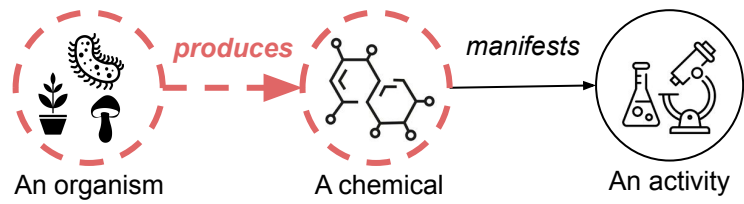
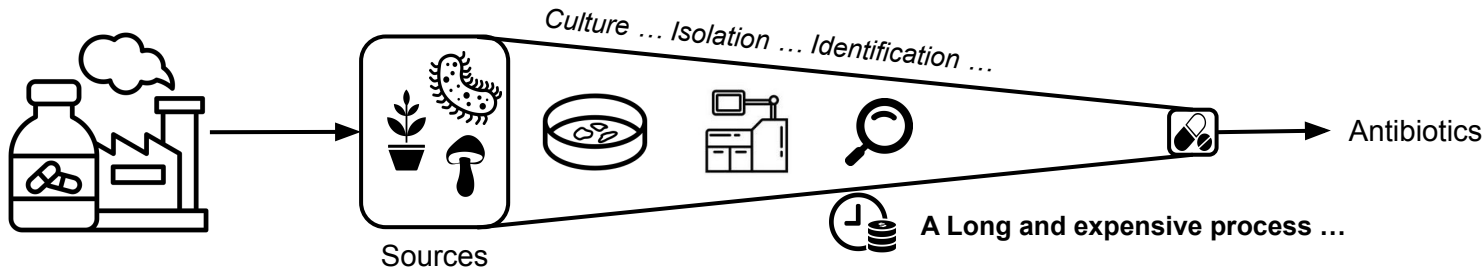
Study context: avoiding rediscoveries



Where is this knowledge ?



Study context: avoiding rediscoveries



Where is this knowledge ?



Scientific literature & patents

LOTUS: An Open Knowledge Base for natural products

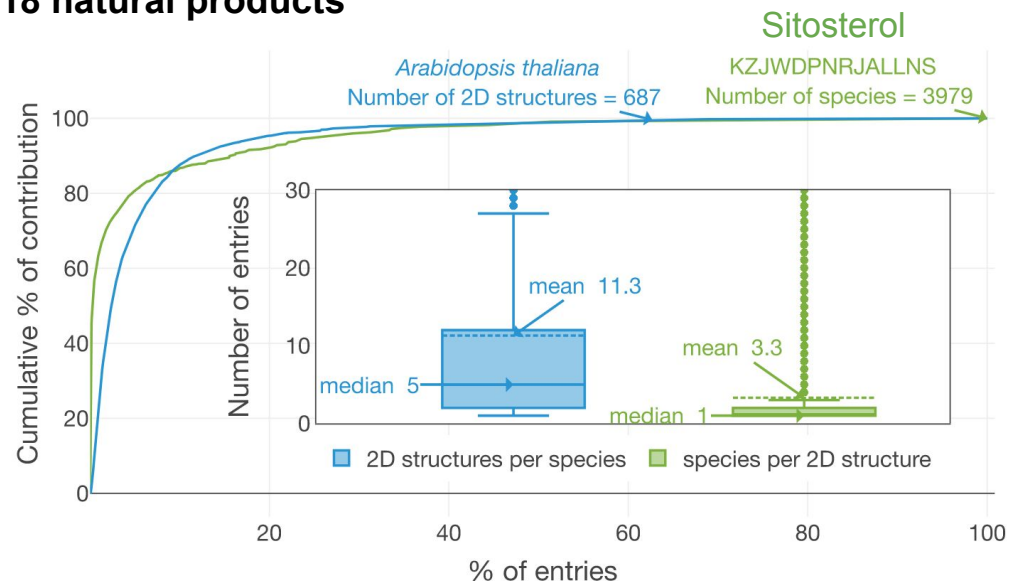
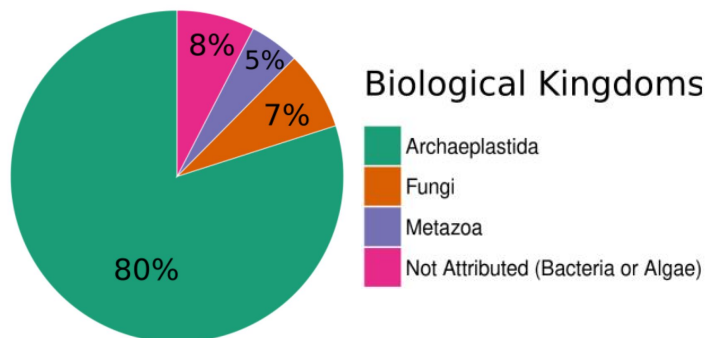


Data on > 276 000 natural products

LOTUS: An Open Knowledge Base for natural products



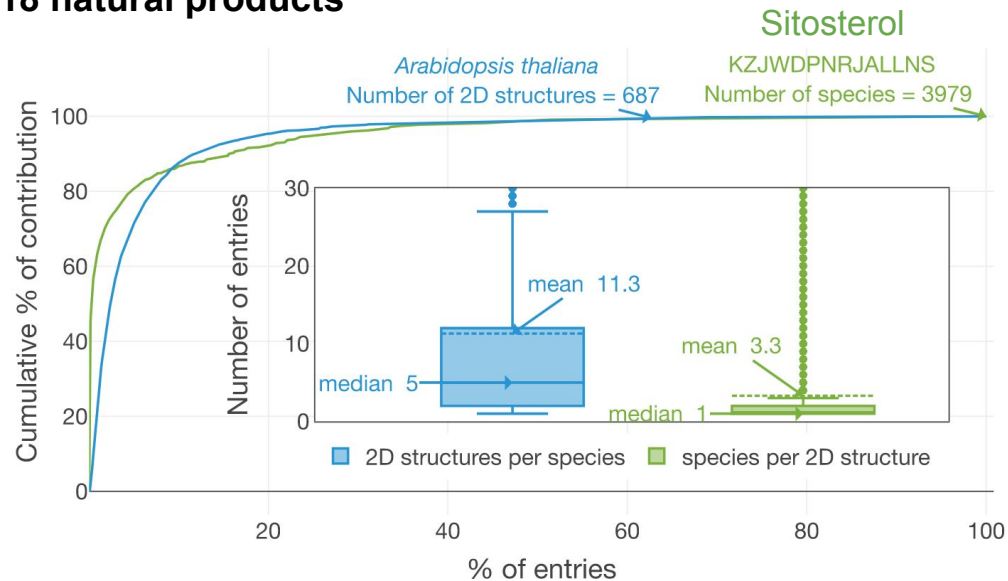
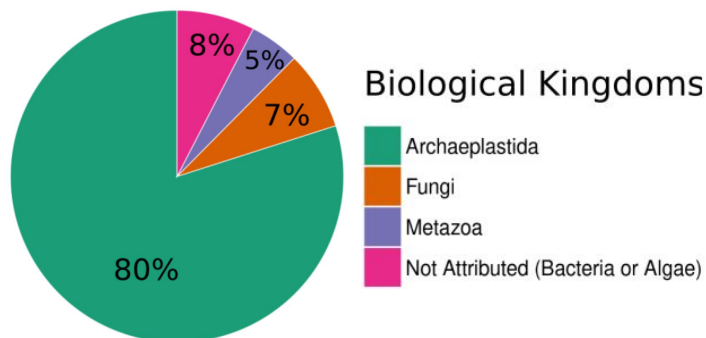
Data on > 276,518 natural products



LOTUS: An Open Knowledge Base for natural products



Data on > 276,518 natural products



Imbalance toward *model organisms*, few data for more *exotic organisms*

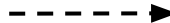
Automatic extraction of relations from the literature



Scientific literature & patents

Abstract

Ochratoxin A (OTA) is a very important mycotoxin, and its research is focused right now on the new findings of OTA, like being a complete carcinogen, information about OTA producers and new exposure sources of OTA. **Citrinin (CIT)** is another important mycotoxin, too, and its research turns towards nephrotoxicity. Both additive and synergistic effects have been described in combination with OTA. OTA is produced in foodstuffs by **Aspergillus** Section Circumdati (**Aspergillus ochraceus**, **A. westerdijkiae**, **A. steynii**) and **Aspergillus** Section Nigri (**Aspergillus carbonarius**, **A. foetidus**, **A. lacticoffeatus**, **A. niger**, **A. sclerotioniger**, **A. tubingensis**), mostly in subtropical and tropical areas. OTA is produced in foodstuffs by **Penicillium verrucosum** and **P. nordicum**, notably in temperate and colder zones. **CIT** is produced in foodstuffs by **Monascus** species (**Monascus purpureus**, **M. ruber**) and **Penicillium** species (**Penicillium citrinum**, **P. expansum**, **P. radicola**, **P. verrucosum**). OTA was frequently found in foodstuffs of both plant origin (e.g., cereal products, coffee, vegetable, liquorice, raisins, wine) and animal origin (e.g., pork/poultry). CIT was also found in foodstuffs of vegetable origin (e.g., cereals, pomaceous fruits, black olive, roasted nuts, spices), food supplements based on rice fermented with red microfungi **Monascus purpureus** and in foodstuffs of animal origin (e.g., cheese).



Extracted relations

Aspergillus ochraceus - Ochratoxin A
Aspergillus westerdijkiae - Ochratoxin A
Aspergillus steynii - Ochratoxin A
...
Monascus purpureus - Citrinin
Penicillium expansum - Citrinin

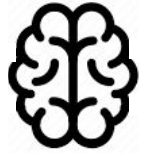
End-to-end
NER / RE



Named Entity Recognition for Bio-entities

How to build a Relation Extraction model - recipe

In the era of machine learning, efficient models are *supervised*



the architecture
(the weights)

How to build a Relation Extraction model - recipe

In the era of machine learning, efficient models are *supervised*



the architecture
(the weights)

+



Ochratoxin A (OTA) is a very important mycotoxin, and its research is focused right now on the new findings of OTA, like being a complete carcinogen, information about OTA producers and new exposure sources of OTA. **Citrinin (CIT)** is another important mycotoxin, too, and its research turns towards nephrotoxicity. Both additive and synergistic effects have been described in combination with OTA. OTA is produced in foodstuffs by **Aspergillus** Section Circumdati (**Aspergillus ochraceus**, **A. westerdijkiae**, **A. steynii**) and **Aspergillus** Section Nigri (**Aspergillus carbonarius**, **A. foetidus**, **A. laticoffeatus**, **A. niger**, **A. sclerotioniger**, **A. tubingensis**), mostly in subtropical and tropical areas. OTA is produced in foodstuffs by **Penicillium verrucosum** and **P. nordicum**, notably in temperate and colder zones. **CIT** is produced in foodstuffs by **Monascus** species (**Monascus purpureus**, **M. ruber**) and **Penicillium** species (**Penicillium citrinum**, **P. expansum**, **P. radicola**, **P. verrucosum**). OTA was frequently found in foodstuffs of both plant origin (e.g., cereal products, coffee, vegetable, liquorice, raisins, wine) and animal origin (e.g., pork/poultry). CIT was also found in foodstuffs of vegetable origin (e.g., cereals, pomaceous fruits, black olive, roasted nuts, spices), food supplements based on rice fermented with red microfungi **Monascus purpureus** and in foodstuffs of animal origin (e.g., cheese).

The data



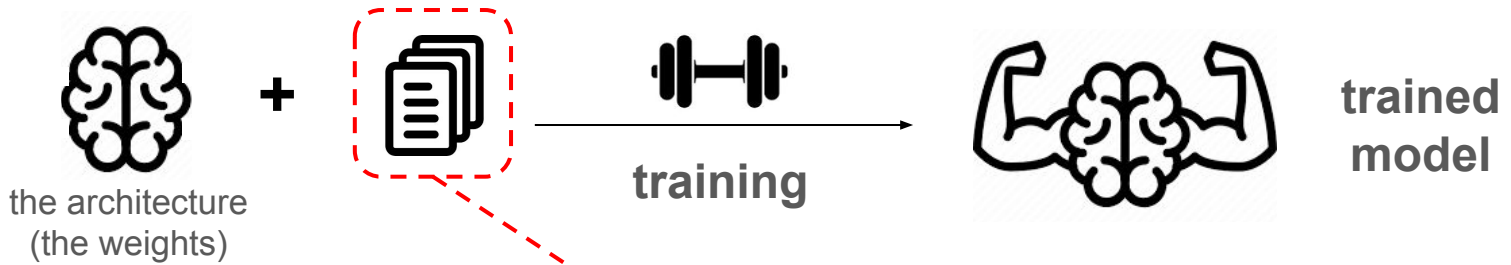
Expected output labels

Aspergillus ochraceus - Ochratoxin A
Aspergillus westerdijkiae - Ochratoxin A
Aspergillus steynii - Ochratoxin A
...
Monascus purpureus - Citrinin
Penicillium expansum - Citrinin

+

How to build a Relation Extraction model - recipe

In the era of machine learning, efficient models are *supervised*



Ochratoxin A (OTA) is a very important mycotoxin, and its research is focused right now on the new findings of OTA, like being a complete carcinogen, information about OTA producers and new exposure sources of OTA. **Citrinin (CIT)** is another important mycotoxin, too, and its research turns towards nephrotoxicity. Both additive and synergistic effects have been described in combination with OTA. OTA is produced in foodstuffs by **Aspergillus** Section Circumdati (**Aspergillus ochraceus**, **A. westerdijkiae**, **A. steynii**) and **Aspergillus** Section Nigri (**Aspergillus carbonarius**, **A. foetidus**, **A. laticoffeatus**, **A. niger**, **A. sclerotioniger**, **A. tubingensis**), mostly in subtropical and tropical areas. OTA is produced in foodstuffs by **Penicillium verrucosum** and **P. nordicum**, notably in temperate and colder zones. **CIT** is produced in foodstuffs by **Monascus** species (**Monascus purpureus**, **M. ruber**) and **Penicillium** species (**Penicillium citrinum**, **P. expansum**, **P. radicola**, **P. verrucosum**). OTA was frequently found in foodstuffs of both plant origin (e.g., cereal products, coffee, vegetable, liquorice, raisins, wine) and animal origin (e.g., pork/poultry). CIT was also found in foodstuffs of vegetable origin (e.g., cereals, pomaceous fruits, black olive, roasted nuts, spices), food supplements based on rice fermented with red microfungi **Monascus purpureus** and in foodstuffs of animal origin (e.g., cheese).

The data 🙏

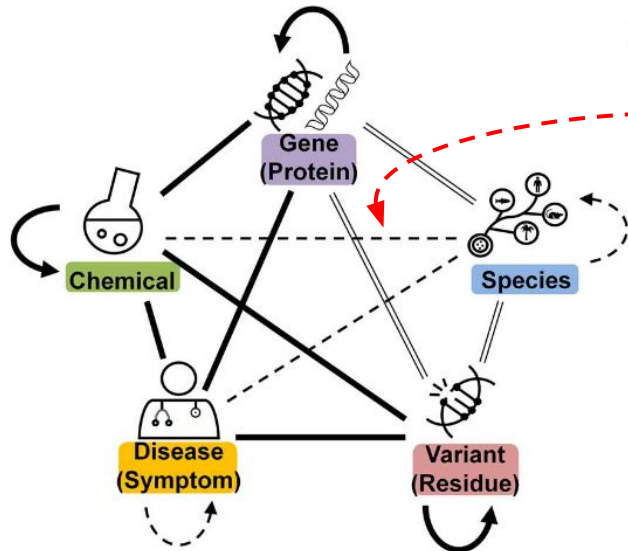
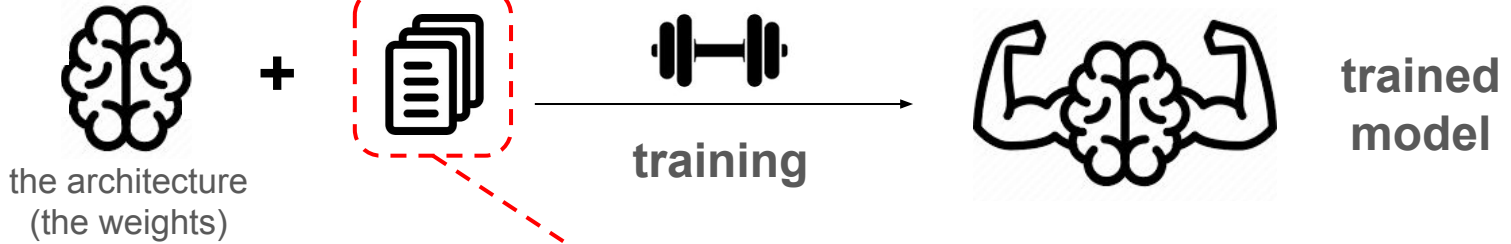
Expected output labels

+

- Aspergillus ochraceus** - Ochratoxin A
- Aspergillus westerdijkiae** - Ochratoxin A
- Aspergillus steynii** - Ochratoxin A
- ...
- Monascus purpureus** - Citrinin
- Penicillium expansum** - Citrinin

How to build a Relation Extraction model - recipe

In the era of machine learning, efficient models are *supervised*



Ochratoxin A (OTA) is a very important mycotoxin, and its research is focused right now on the new findings of OTA, like being a complete carcinogen, information about OTA producers and new exposure sources of OTA. **Citrinin (CIT)** is another important mycotoxin, too, and its research turns towards nephrotoxicity. Both additive and synergistic effects have been described in combination with OTA. OTA is produced in foodstuffs by *Aspergillus* Section Circumdati (*Aspergillus ochraceus*, *A. westerdijkiae*, *A. steynii*) and *Aspergillus* Section Nigri (*Aspergillus carbonarius*, *A. foetidus*, *A. laticoffeatus*, *A. niger*, *A. sclerotioniger*, *A. tubingensis*), mostly in subtropical and tropical areas. OTA is produced in foodstuffs by *Penicillium verrucosum* and *P. nordicum*, notably in temperate and colder zones. **CIT** is produced in foodstuffs by *Monascus* species (*Monascus purpureus*, *M. ruber*) and *Penicillium* species (*Penicillium citrinum*, *P. expansum*, *P. radicicola*, *P. verrucosum*). OTA was frequently found in foodstuffs of both plant origin (e.g., cereal products, coffee, vegetable, liquorice, raisins, wine) and animal origin (e.g., pork/poultry). CIT was also found in foodstuffs of vegetable origin (e.g., cereals, pomaceous fruits, black olive, roasted nuts, spices), food supplements based on rice fermented with red microfungi *Monascus purpureus* and in foodstuffs of animal origin (e.g., cheese).

The data 🙏

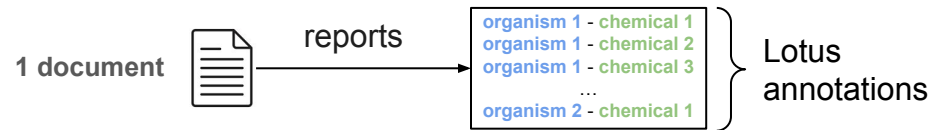
Expected output labels

Aspergillus ochraceus - Ochratoxin A
Aspergillus westerdijkiae - Ochratoxin A
Aspergillus steynii - Ochratoxin A
...
Monascus purpureus - Citrinin
Penicillium expansum - Citrinin

But there is no curated datasets ...

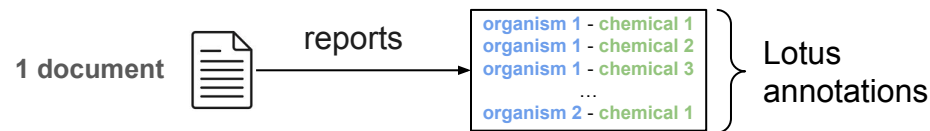
But what about LOTUS ?

LOTUS  Harmonization / Processing / Validation / Dissemination



But what about LOTUS ?

LOTUS Harmonization / Processing / Validation / Dissemination



5-hydroxytryptamine-derived alkaloids from two marine sponges of the genus Hyrtios.

Salmoun M, Devijver C ... van Soest RW • J. Nat. Prod.

[Add to Collection](#) [BiocXML](#)

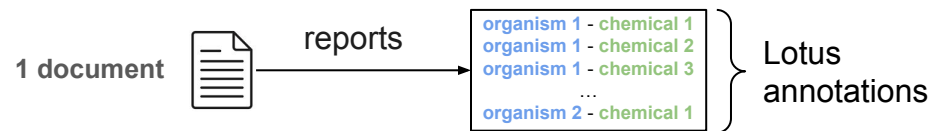
Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties. *H. erectus* contained hyrtiosulawesine (4), a new beta-carboline alkaloid, together with the already known alkaloids 5-hydroxyindole-3-carbaldehyde (1), hyrtiosin B (2), and 5-hydroxy-3-(2-hydroxyethyl)indole (3). *H. reticulatus* contained the novel derivative 1,6-dihydroxy-1,2,3,4-tetrahydro-beta-carboline (11) together with serotonin (5), 6-hydroxy-1-methyl-1,2,3,4-tetrahydro-beta-carboline (7), and 6-hydroxy-3,4-dihydro-1-oxo-beta-carboline (9).

- Hyrtios erectus* - Hyrtiosulawesine ✓
- Hyrtios erectus* - 5-hydroxy-1H-indole-3-carbaldehyde ✗
- Hyrtios erectus* - 1,2-bis(5-hydroxy-1H-indol-3-yl)ethane-1,2-dione ✗
- Hyrtios erectus* - 5-Hydroxytryptophol ✗
- Hyrtios reticulatus* - (1R)-2,3,4,9-tetrahydro-1H-pyrido[3,4-b]indole-1,6-diol ✗
- Hyrtios reticulatus* - Serotonin ✓
- Hyrtios reticulatus* - (1R)-1-methyl-2,3,4,9-tetrahydro-1H-pyrido[3,4-b]indol-6-ol ✗
- Hyrtios reticulatus* - 2,3,4,9-Tetrahydro-6-hydroxy-1H-pyrido(3,4-b)indol-1-one ✗
- Hyrtios reticulatus* - (S)-6-Hydroxytetrahydroharman ✗

Discrepancies between text and expected output labels

But what about LOTUS ?

LOTUS Harmonization / Processing / Validation / Dissemination



5-hydroxytryptamine-derived alkaloids from two marine sponges of the genus Hyrtios.

Salmoun M, Devijver C ... van Soest RW • J. Nat. Prod.

[Add to Collection](#) [BiocXML](#)

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties. *H. erectus* contained hyrtiosulawesine (4), a new beta-carboline alkaloid, together with the already known alkaloids 5-hydroxyindole-3-carbaldehyde (1), hyrtiosin B (2), and 5-hydroxy-3-(2-hydroxyethyl)indole (3). *H. reticulatus* contained the novel derivative 1,6-dihydroxy-1,2,3,4-tetrahydro-beta-carboline (11) together with serotonin (5), 6-hydroxy-1-methyl-1,2,3,4-tetrahydro-beta-carboline (7), and 6-hydroxy-3,4-dihydro-1-oxo-beta-carboline (9).

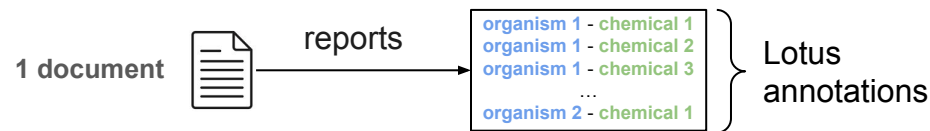
- Hyrtios erectus* - Hyrtiosulawesine ✓
- Hyrtios erectus* - 5-hydroxy-1H-indole-3-carbaldehyde ✗
- Hyrtios erectus* - 1,2-bis(5-hydroxy-1H-indol-3-yl)ethane-1,2-dione ✗
- Hyrtios erectus* - 5-Hydroxytryptophol ✗
- Hyrtios reticulatus* - (1R)-2,3,4,9-tetrahydro-1H-pyrido[3,4-b]indole-1,6-diol ✗
- Hyrtios reticulatus* - Serotonin ✓
- Hyrtios reticulatus* - (1R)-1-methyl-2,3,4,9-tetrahydro-1H-pyrido[3,4-b]indol-6-ol ✗
- Hyrtios reticulatus* - 2,3,4,9-Tetrahydro-6-hydroxy-1H-pyrido(3,4-b)indol-1-one ✗
- Hyrtios reticulatus* - (S)-6-Hydroxytetrahydroharman ✗

Discrepancies between text and expected output labels

Not design for building NLP related datasets

But what about LOTUS ?

LOTUS Harmonization / Processing / Validation / Dissemination



5-hydroxytryptamine-derived alkaloids from two marine sponges of the genus Hyrtios.

Salmoun M, Devijver C ... van Soest RW • J. Nat. Prod.

[Add to Collection](#) [BiocXML](#)

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties. *H. erectus* contained hyrtiosulawesine (4), a new beta-carboline alkaloid, together with the already known alkaloids 5-hydroxyindole-3-carbaldehyde (1), hyrtiosin B (2), and 5-hydroxy-3-(2-hydroxyethyl)indole (3). *H. reticulatus* contained the novel derivative 1,6-dihydroxy-1,2,3,4-tetrahydro-beta-carboline (11) together with serotonin (5), 6-hydroxy-1-methyl-1,2,3,4-tetrahydro-beta-carboline (7), and 6-hydroxy-3,4-dihydro-1-oxo-beta-carboline (9).

- Hyrtios erectus* - Hyrtiosulawesine ✓
- Hyrtios erectus* - 5-hydroxy-1H-indole-3-carbaldehyde ✗
- Hyrtios erectus* - 1,2-bis(5-hydroxy-1H-indol-3-yl)ethane-1,2-dione ✗
- Hyrtios erectus* - 5-Hydroxytryptophol ✗
- Hyrtios reticulatus* - (1R)-2,3,4,9-tetrahydro-1H-pyrido[3,4-b]indole-1,6-diol ✗
- Hyrtios reticulatus* - Serotonin ✓
- Hyrtios reticulatus* - (1R)-1-methyl-2,3,4,9-tetrahydro-1H-pyrido[3,4-b]indol-6-ol ✗
- Hyrtios reticulatus* - 2,3,4,9-Tetrahydro-6-hydroxy-1H-pyrido(3,4-b)indol-1-one ✗
- Hyrtios reticulatus* - (S)-6-Hydroxytetrahydroharman ✗

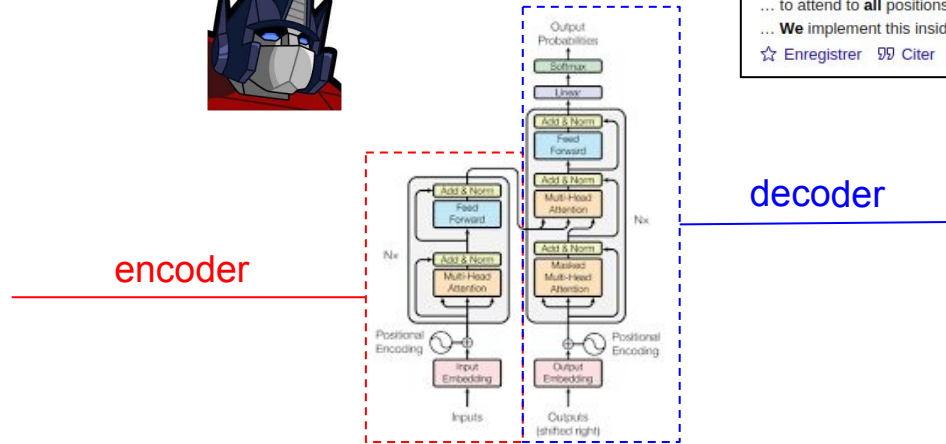
Discrepancies between text and expected output labels

Not design for building NLP related datasets

Can LLM help us in this context ?

What are Large Language Models ?

Transformers



Attention is all you need

[A Vaswani, N Shazeer, N Parmar...](#) - Advances in neural ..., 2017 - proceedings.neurips.cc

... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent

... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...

☆ Enregistrer Citer Cité 110783 fois Autres articles Les 62 versions »»

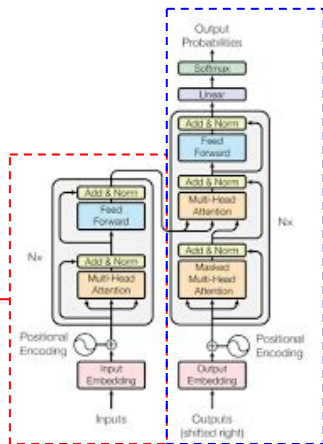
What are Large Language Models ?

Transformers



Representation

encoder



Attention is all you need

[A Vaswani, N Shazeer, N Parmar...](#) - Advances in neural ..., 2017 - proceedings.neurips.cc

... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent ... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...

☆ Enregistrer Citer Cité 110783 fois Autres articles Les 62 versions »»

decoder

Bert: Pre-training of deep **bidirectional** transformers for language understanding

J De
... B
K Clark, M
2018
... back-p
☆ E
... gains f

Electra: Pre-training text encoders as discriminators rather than generators

Roberta: A robustly optimized bert pretraining approach

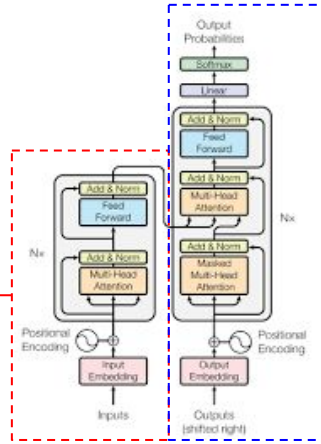
[Y Liu, M Ott, N Goyal, J Du, M Joshi, D Chen...](#) - arXiv preprint arXiv ..., 2019 - arxiv.org

... We find that **BERT** was significantly undertrained and propose an improved ... **BERT** models, which we call **RoBERTa**, that can match or exceed the performance of all of the post-**BERT** ...

☆ Enregistrer Citer Cité 9454 fois Autres articles Les 6 versions »»

What are Large Language Models ?

Transformers



Representation

encoder



Attention is all you need

[A Vaswani, N Shazeer, N Parmar...](#) - Advances in neural ..., 2017 - proceedings.neurips.cc

... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent

... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...

☆ Enregistrer 📄 Citer Cité 110783 fois Autres articles Les 62 versions 📄

decoder

• Generation

Generative Pre-training Transformer

GPT Series Models



Bert: Pre-training of deep **bidirectional** transformers for language understanding

Electra: Pre-training text encoders as discriminators rather than generators

Roberta: A robustly optimized bert pretraining approach

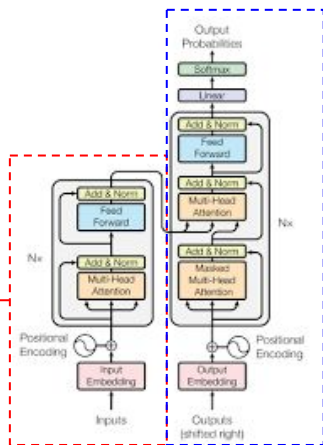
[Y Liu, M Ott, N Goyal, J Du, M Joshi, D Chen...](#) - arXiv preprint arXiv ..., 2019 - arxiv.org

... We find that **BERT** was significantly undertrained and propose an improved ... **BERT** models, which we call **RoBERTa**, that can match or exceed the performance of all of the post-BERT ...

☆ Enregistrer 📄 Citer Cité 9454 fois Autres articles Les 6 versions 📄

What are Large Language Models ?

Transformers



Representation

encoder



T5

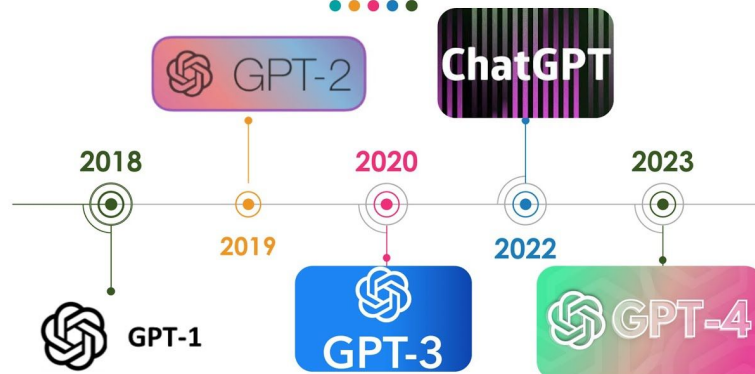
Attention is all you need

[A Vaswani, N Shazeer, N Parmar...](#) - Advances in neural ..., 2017 - proceedings.neurips.cc
... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent ... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...
☆ Enregistrer Citer Cité 110783 fois Autres articles Les 62 versions »

decoder • Generation

Generative Pre-training Transformer

GPT Series Models



Bert: Pre-training of deep **bidirectional** transformers for language understanding

Electra: Pre-training text encoders as discriminators rather than generators

Roberta: A robustly optimized bert pretraining approach

[Y Liu, M Ott, N Goyal, J Du, M Joshi, D Chen...](#) - arXiv preprint arXiv ..., 2019 - arxiv.org

... We find that **BERT** was significantly undertrained and propose an improved ... **BERT** models, which we call **RoBERTa**, that can match or exceed the performance of all of the post-BERT ...

☆ Enregistrer Citer Cité 9454 fois Autres articles Les 6 versions »

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



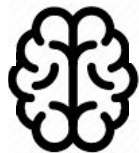
the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



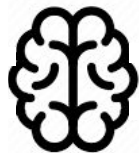
the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



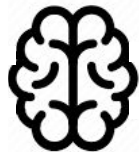
the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



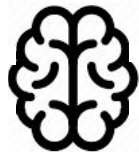
the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



the architecture
(the weights - Billions)

Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral

Pre-trained / Foundation models

GPT: Generative **Pre-trained** Models



Task: Predict the next word

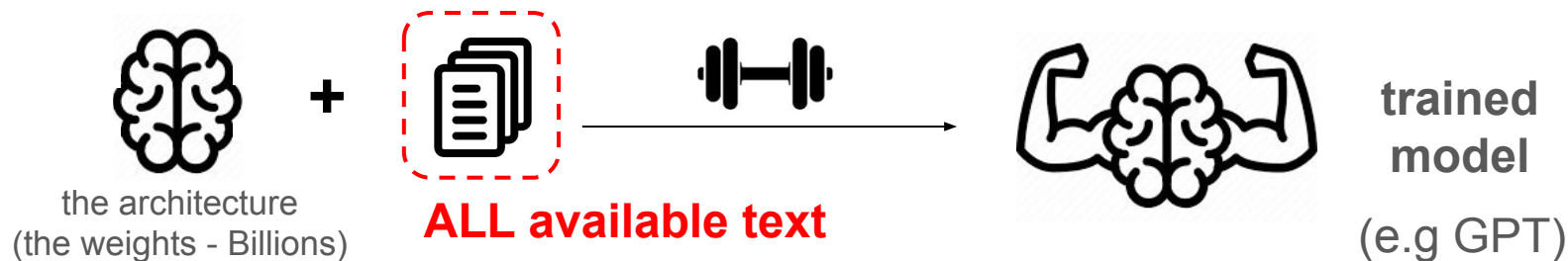
Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties



Learn a **representation** of the text

Pre-trained / Foundation models

GPT: Generative Pre-trained Models



Task: Predict the next word

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties



Learn a **representation** of the text

How Can We Know What Language Models Know?

Zhengbao Jiang^{1*} Frank F. Xu^{1*} Jun Araki² Graham Neubig¹

Language Technologies Institute, Carnegie Mellon University¹

Bosch Research North America²

{zhengba j, fangzhex, gneubig}@cs.cmu.edu jun.araki@us.bosch.com

Evaluating Open-Domain Question Answering in the Era of Large Language Models

Ehsan Kamaloo[♦] Nouha Dziri[♣] Charles L. A. Clarke[♣] Davood Rafiei[♦]

Knowledgeable or Educated Guess? Revisiting Language Models as Knowledge Bases

Boxi Cao^{1,3}, Hongyu Lin¹, Xianpei Han^{1,2}, Le Sun^{1,2}

Lingyong Yan^{1,3}, Meng Liao¹, Tong Xue¹, Jin Xu⁴

¹Chinese Information Processing Laboratory ²State Key Laboratory of Computer Science
Institute of Software, Chinese Academy of Sciences, Beijing, China

³University of Chinese Academy of Sciences, Beijing, China

⁴Data Quality Team, WeChat, Tencent Inc., China

{boxi2020, hongyu, xianpei, sunle, lingyong2014}@iscas.ac.cn

{maricoliao, xavierxue, jinxxu}@tencent.com

Text as a projection of the world: real knowledge ?

Evaluating Open-QA Evaluation

Cunxiang Wang¹, Sirui Cheng², Qipeng Guo³, Yuanhao Yue⁴, Bowen Ding¹,
Zhikun Xu⁴, Yidong Wang¹, Xiangkun Hu³, Zheng Zhang³, and Yue Zhang^{1†}

¹School of Engineering, Westlake University, China

²Northeastern University, China; ³Amazon AWS AI; ⁴Fudan University, China

{wangcunxiang, zhanyue}@westlake.edu.cn

Statistical Knowledge Assessment for Large Language Models

Qingxiu Dong¹, Jingjing Xu², Lingpeng Kong³, Zhifang Sui¹ and Lei Li⁴

¹National Key Laboratory for Multimedia Information Processing,

School of Computer Science, Peking University

²Shanghai AI Lab ³The University of Hong Kong ⁴Carnegie Mellon University
dqx@stu.pku.edu.cn, {jingjingxu, szf}@pku.edu.cn, lpk@cs.hku.hk, leili@cs.cmu.edu

The Internal State of an LLM Knows When It's Lying

Amos Azaria
School of Computer Science,
Ariel University, Israel

Tom Mitchell
Machine Learning Dept.,
Carnegie Mellon University, Pittsburgh, PA

Measuring and Modifying Factual Knowledge in Large Language Models

Pouya Pezeshkpour
Megagon Labs
pouya@megagon.ai



KoLA: Carefully Benchmarking World Knowledge of Large Language Models

Jifan Yu², Xiaozhi Wang², Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv,
Hao Peng, Zijun Yao, Xiaohan Zhang, Hanming Li, Chunyang Li,
Zheyuan Zhang, Yushi Bai, Yantao Liu, Amy Xin, Nianyi Lin, Kaifeng Yun,
Linlu Gong, Jianhui Chen, Zhili Wu, Yunjia Qi, Weikai Li, Yong Guan,
Kaisheng Zeng, Ji Qi, Hailong Jin, Jinxin Liu, Yu Gu, Yuan Yao, Ning Ding,
Lei Hou, Zhiyuan Liu, Bin Xu, Jie Tang, Juanzi Li¹

¹Tsinghua University
Beijing, China 100084

kola-benchmark@googlelegroups.com

Evaluating Open-QA Evaluation

Cunxiang Wang¹, Sirui Cheng², Qipeng Guo³, Yuanhao Yue⁴, Bowen Ding¹,
Zhikun Xu⁴, Yidong Wang¹, Xiangkun Hu³, Zheng Zhang³, and Yue Zhang^{1†}

¹School of Engineering, Westlake University, China

²Northeastern University, China; ³Amazon AWS AI; ⁴Fudan University, China
{wangcunxiang, zhangyue}@westlake.edu.cn

Language Models as Knowledge Bases?

Fabio Petroni¹ Tim Rocktäschel^{1,2} Patrick Lewis^{1,2} Anton Bakhtin¹
Yuxiang Wu^{1,2} Alexander H. Miller¹ Sebastian Riedel^{1,2}

¹Facebook AI Research

²University College London

{fabiopetroni, rockt, plewis, yolo, yuxiangwu, ahm, sriedel}@fb.com

LLM as Knowledge bases ?

 **> 35 million citations**  *Hard to process ...*

*But, this literature has already been “**digested**” during pre-training !*

LLM as Knowledge bases ?



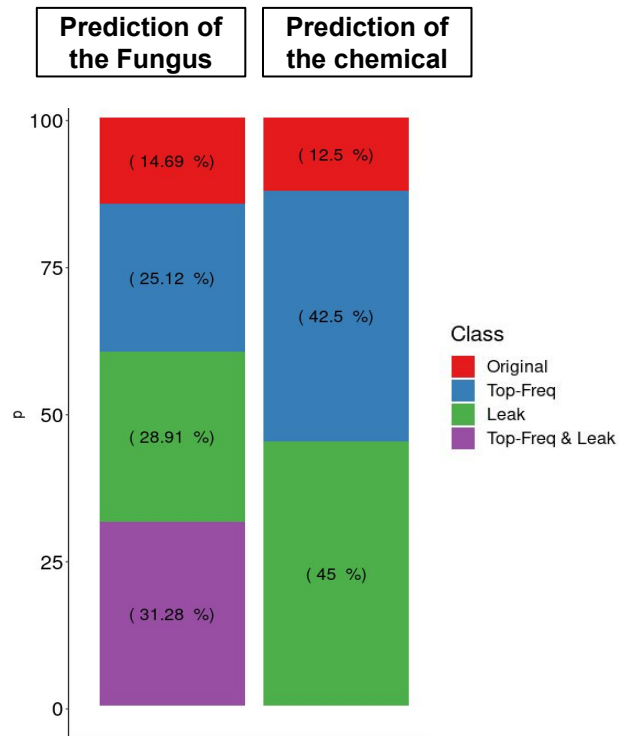
> 35 million citations



Hard to process ...

But, this literature has already been “**digested**” during pre-training !

- Short answer: **No ...**
- **Low global retrieval performances**
- **Correct predictions are mostly leaks**
E.g: *Monascus Pilosus produces monascin*
- **Some predictions are generics**
E.g: *Aspergillus Niger* or *Ergosterol*



LLM as Knowledge bases ?



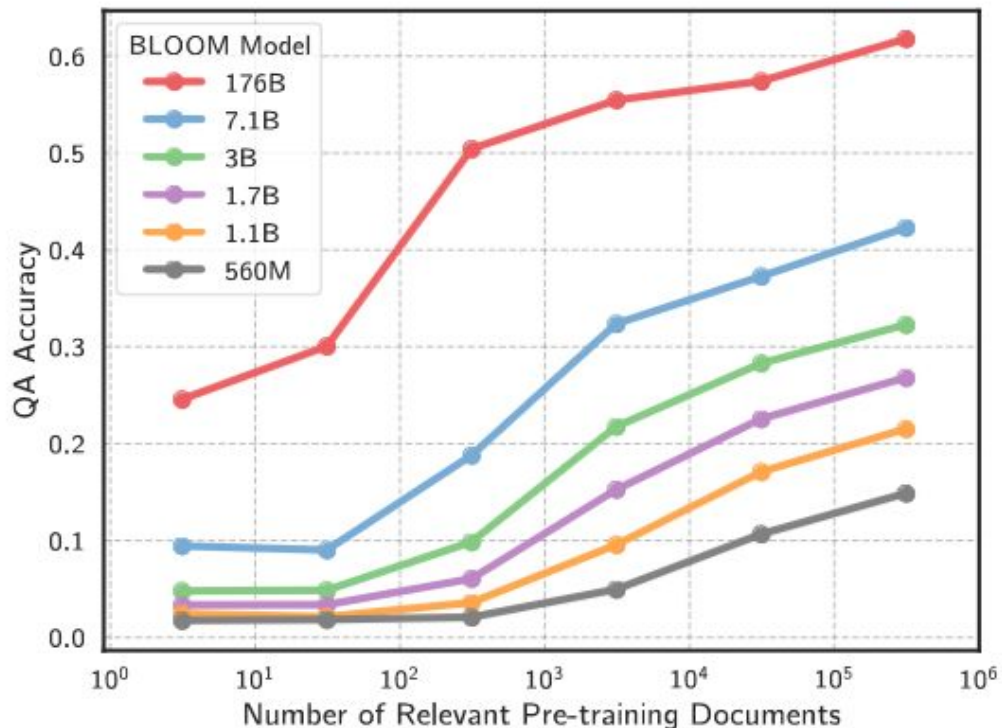
> 35 million citations



Hard to process ...

But, this literature has already been “**digested**” during pre-training !

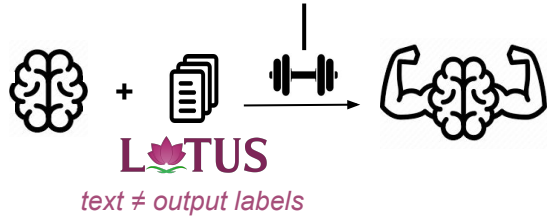
- Short answer: **No ...**
- **Low global retrieval performances**
- **Correct predictions are mostly leaks**
E.g: *Monascus Pilosus produces monascin*
- **Some predictions are generics**
E.g: *Aspergillus Niger* or *Ergosterol*



We need to extract the relation from the text !

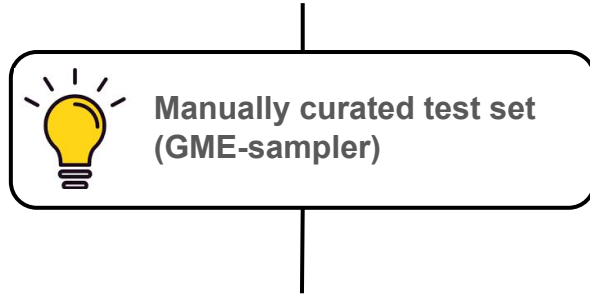
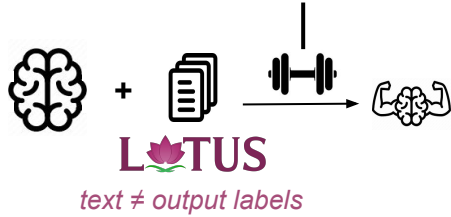
We need to extract the relation from the text !

Supervised



We need to extract the relation from the text !

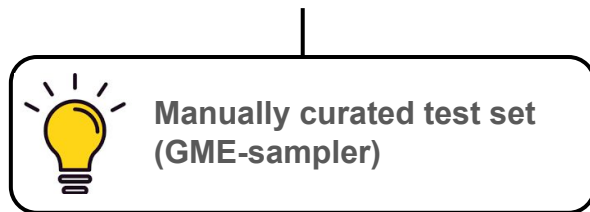
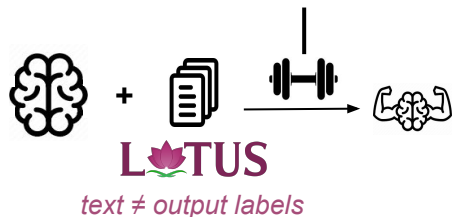
Supervised



	precision	recall	f1
Seq2rel	47.3	5.8	10.4
GPT-2	44.8	21.7	29.3
BioGPT	42.2	26.5	32.5

We need to extract the relation from the text !

Supervised



	precision	recall	f1
Seq2rel	47.3	5.8	10.4
GPT-2	44.8	21.7	29.3
BioGPT	42.2	26.5	32.5

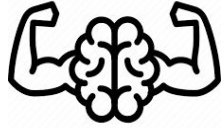
(Weakly) supervised

Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*	
Jared Kaplan†	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Mateusz Litwin	Scott Gray
	Benjamin Chess	Jack Clark	Christopher Berner	
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei	
				OpenAI

“The model is conditioned on a natural language instruction and/or a few demonstrations of the task and is then expected to complete further instances of the task simply by predicting what comes next.”

LLM are Few shot learners



Use the representation learned during pre-training

Poor English input: I eated the purple berries.

Good English output: I ate the purple berries.

Poor English input: Thank you for picking me as your designer. I'd appreciate it.

Good English output: Thank you for choosing me as your designer. I appreciate it.

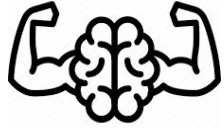
Poor English input: The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

Good English output: The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

Poor English input: I'd be more than happy to work with you in another project.

Good English output: LLM completing ...

LLM are Few shot learners



Use the representation learned during pre-training

Poor English input: I eated the purple berries.

Good English output: I ate the purple berries.

Poor English input: Thank you for picking me as your designer. I'd appreciate it.

Good English output: Thank you for choosing me as your designer. I appreciate it.

Poor English input: The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

Good English output: The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

Poor English input: I'd be more than happy to work with you in another project.

Good English output: I'd be more than happy to work with you on another project.

LLM are Few shot learners

Instruction

The task is to extract relations between organisms and chemicals from the input text.

INPUT: The antimicrobially active EtOH extracts of *Maytenus heterophylla* yielded a new dihydroagarofuran alkaloid, *1beta-acetoxy-9alpha-benzoyloxy-dihydroagarofuran*, together with the known compounds *beta-amyrin*, *maytenfolic acid*, ...

OUTPUT: *Maytenus heterophylla* produces *1beta-acetoxy-9alpha-benzoyloxy-dihydroagarofuran*. *Maytenus heterophylla* produces *beta-amyrin*. *Maytenus heterophylla* produces *maytenfolic acid*.

INPUT: Ten new ergosteroids, *gloeophyllins A-J* (1-10), have been isolated from the solid cultures of *Gloeophyllum abietinum*.

OUTPUT: *Gloeophyllum abietinum* produces *gloeophyllin A*. *Gloeophyllum abietinum* produces *gloeophyllin B*. *Gloeophyllum abietinum* produces *gloeophyllin C*. *Gloeophyllum abietinum* produces *gloeophyllin D*. *Gloeophyllum abietinum* produces *gloeophyllin I*. *Gloeophyllum abietinum* produces *gloeophyllin J*.

INPUT: The present work describes the isolation of the cyclic peptides *geodiamolides A, B, H and I* (1-4) from *G. corticostylifera* and their anti-proliferative effects against sea urchin eggs and human breast cancer cell lineages.

OUTPUT: *G. corticostylifera* produces *geodiamolide A*. *G. corticostylifera* produces *geodiamolide B* [...]

INPUT: Four new cyclic peptides, *patellamide G* (2) and *ulithiacyclamides E-G* (3-5), along with the known *patellamides A-C* (6-8) and *ulithiacyclamide B* (9), were isolated from the ascidian *Lissoclinum patella* collected in Pohnpei, Federated States of Micronesia.

OUTPUT: *Lissoclinum patella* produces *patellamide G*. *Lissoclinum patella* produces *ulithiacyclamide E*. *Lissoclinum patella* produces *ulithiacyclamide F*. *Lissoclinum patella* produces *ulithiacyclamide B*.

INPUT: Chemical investigation of *Trogopterus faeces* has led to the isolation of seven *flavonoids*. Their structures were elucidated by chemical and spectral analyses. In an anticoagulative assay, three *kaempferol coumaroyl rhamnosides* had significant antithrombin activity. This is the first report on the occurrence of *flavonoid glycosides* in *Trogopterus faeces*.

OUTPUT: *Trogopterus faeces* produces *flavonoids*. *Trogopterus faeces* produces *kaempferol coumaroyl rhamnosides*. *Trogopterus faeces* produces *flavonoid glycosides*.

Demonstrations

↓
Archetypal sentences

New Instance → **INPUT: ** Abstract ****

To fill → **OUTPUT: [LLM completing ...]**

LLM are Few shot learners

Instruction

The task is to extract relations between organisms and chemicals from the input text.

INPUT: The antimicrobially active EtOH extracts of *Maytenus heterophylla* yielded a new dihydroagarofuran alkaloid, **1beta-acetoxy-9alpha-benzoyloxy-dihydroagarofuran**, together with the known compounds **beta-amyrin**, **maytenfolic acid**, ...

OUTPUT: *Maytenus heterophylla* produces **1beta-acetoxy-9alpha-benzoyloxy-dihydroagarofuran**. *Maytenus heterophylla* produces **beta-amyrin**. *Maytenus heterophylla* produces **maytenfolic acid**.

INPUT: Ten new ergosteroids, **gloeophyllins A-J** (1-10), have been isolated from the solid cultures of *Gloeophyllum abietinum*.

OUTPUT: *Gloeophyllum abietinum* produces **gloeophyllin A**. *Gloeophyllum abietinum* produces **gloeophyllin B**. *Gloeophyllum abietinum* produces **gloeophyllin C**. *Gloeophyllum abietinum* produces **gloeophyllin D**. *Gloeophyllum abietinum* produces **gloeophyllin I**. *Gloeophyllum abietinum* produces **gloeophyllin J**.

INPUT: The present work describes the isolation of the cyclic peptides **geodiamolides A, B, H and I** (1-4) from *G. corticostylifera* and their anti-proliferative effects against sea urchin eggs and human breast cancer cell lineages.

OUTPUT: *G. corticostylifera* produces **geodiamolide A**. *G. corticostylifera* produces **geodiamolide B** [...]

INPUT: Four new cyclic peptides, **patellamide G** (2) and **ulithiacyclamides E-G** (3-5), along with the known **patellamides A-C** (6-8) and **ulithiacyclamide B** (9), were isolated from the ascidian *Lissoclinum patella* collected in Pohnpei, Federated States of Micronesia.

OUTPUT: *Lissoclinum patella* produces **patellamide G**. *Lissoclinum patella* produces **ulithiacyclamide E**. *Lissoclinum patella* produces **ulithiacyclamide F**. *Lissoclinum patella* produces **ulithiacyclamide B**.

INPUT: Chemical investigation of *Trogopterus faeces* has led to the isolation of seven **flavonoids**. Their structures were elucidated by chemical and spectral analyses. In an anticoagulative assay, three **kaempferol coumaroyl rhamnosides** had significant antithrombin activity. This is the first report on the occurrence of **flavonoid glycosides** in *Trogopterus faeces*.

OUTPUT: *Trogopterus faeces* produces **flavonoids**. *Trogopterus faeces* produces **kaempferol coumaroyl rhamnosides**. *Trogopterus faeces* produces **flavonoid glycosides**.

New Instance

INPUT: Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties. *H. erectus* contained **hyrtiosulawesine** (4), a new beta-carboline alkaloid, together with the already known alkaloids **5-hydroxyindole-3-carbaldehyde** (1), **hyrtiosin B** (2), and **5-hydroxy-3-(2-hydroxyethyl)indole** (3). *H. reticulatus* contained the novel derivative **1,6-dihydroxy-1,2,3,4-tetrahydro-beta-carboline** (11) together with **serotonin** (5), **6-hydroxy-1-methyl-1,2,3,4-tetrahydro-beta-carboline** (7), and **6-hydroxy-3,4-dihydro-1-oxo-beta-carboline** (9).

To fill

OUTPUT: *Hyrtios erectus* produces **hyrtiosulawesine**, *Hyrtios erectus* produces **5-hydroxyindole-3-carbaldehyde**. *Hyrtios erectus* produces **hyrtiosin B**

Demonstrations

↓
Archetypal sentences


We need to extract the relation from the text !

Supervised



LOTUS

text ≠ output labels

 Manually curated test set (GME-sampler)

	precision	recall	f1
Seq2rel	47.3	5.8	10.4
GPT-2	44.8	21.7	29.3
BioGPT	42.2	26.5	32.5

(Weakly) supervised

Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*	
Jared Kaplan*	Pranav Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariad Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Mateusz Litwin	Scott Gray
	Benjamin Chess	Jack Clark	Christopher Berner	
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei	

OpenAI

 Only Open Models !

Competitive
↔
with only 5 examples

LLM	precision	recall	f1
Llama-7B	27.0	9.04	13.55
Llama-13B	35.64	23.64	28.49
Llama-30B	38.51	<u>23.24</u>	28.99
Llama-65B	<u>40.16</u>	22.97	<u>29.23</u>
Alpaca-7B	15.14	2.21	5.86
Vicuna-13B	38.4	20.43	26.48

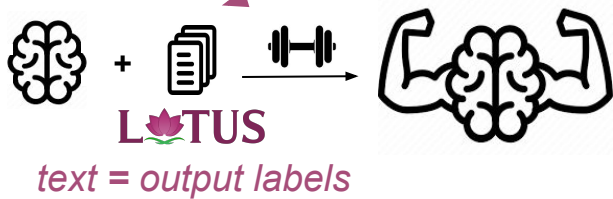
The data is the main bottleneck



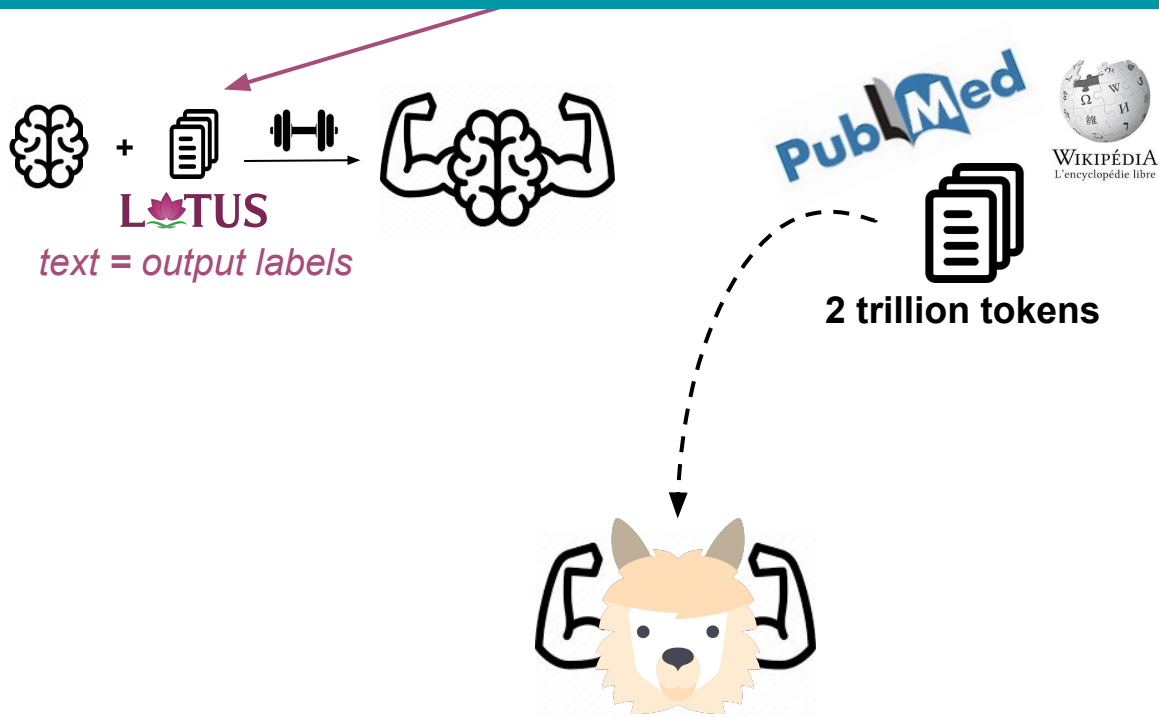
LOTUS

text \neq output labels

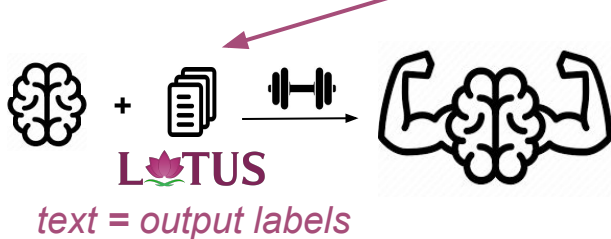
The data is the main bottleneck



The data is the main bottleneck



The data is the main bottleneck



input text

Ochratoxin A (OTA) is a very important mycotoxin, and its research is focused right now on the new findings of OTA, like being a complete carcinogen, information about OTA producers and new exposure sources of OTA. **Citrinin (CIT)** is another important mycotoxin, too, and its research turns towards nephrotoxicity. Both additive and synergistic effects have been described in combination with OTA. OTA is produced in foodstuffs by **Aspergillus** Section Circumdati (**Aspergillus ochraceus**, **A. westerdijkiae**, **A. steynii**) and **Aspergillus** Section Nigri (**Aspergillus carbonarius**, **A. foetidus**, **A. lacticofeatus**, **A. niger**, **A. sclerotiumiger**, **A. tubingensis**), mostly in subtropical and tropical areas. OTA is produced in foodstuffs by **Penicillium verrucosum** and **P. nordicum**, notably in temperate and colder zones. **CIT** is produced in foodstuffs by **Monascus** species (**Monascus purpureus**, **M. ruber**) and **Penicillium** species (**Penicillium citrinum**, **P. expansum**, **P. radicicola**, **P. verrucosum**). OTA was frequently found in foodstuffs of both plant origin (e.g., cereal products, coffee, vegetable, liquorice, raisins, wine) and animal origin (e.g., pork/poultry). CIT was also found in foodstuffs of vegetable origin (e.g., cereals, pomaceous fruits, black olive, roasted nuts, spices), food supplements based on rice fermented with red microfungi **Monascus purpureus** and in foodstuffs of animal origin (e.g., cheese).



2 trillion tokens

Expected relations

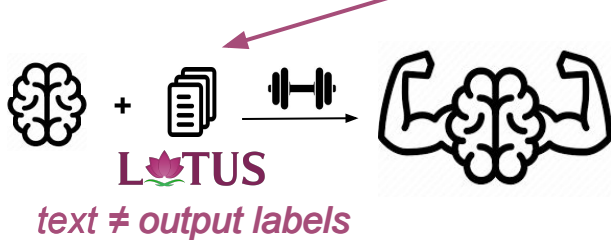
task



hard

Aspergillus ochraceus - **Ochratoxin A**
Aspergillus westerdijkiae - **Ochratoxin A**
Aspergillus steynii - **Ochratoxin A**
...
Monascus purpureus - **Citrinin**
Penicillium expansum - **Citrinin**

The data is the main bottleneck



input text



Expected relations

Ochratoxin A (OTA) is a very important mycotoxin, and its research is focused right now on the new findings of OTA, like being a complete carcinogen, information about OTA producers and new exposure sources of OTA. **Citrinin (CIT)** is another important mycotoxin, too, and its research turns towards nephrotoxicity. Both additive and synergistic effects have been described in combination with OTA. OTA is produced in foodstuffs by *Aspergillus* Section Circumdati (*Aspergillus ochraceus*, *A. westerdijkiae*, *A. steynii*) and *Aspergillus* Section Nigri (*Aspergillus carbonarius*, *A. foetidus*, *A. lacticofeatus*, *A. niger*, *A. sclerotiorum*, *A. tubingensis*), mostly in subtropical and tropical areas. OTA is produced in foodstuffs by *Penicillium verrucosum* and *P. nordicum*, notably in temperate and colder zones. **CIT** is produced in foodstuffs by *Monascus* species (*Monascus purpureus*, *M. ruber*) and *Penicillium* species (*Penicillium citrinum*, *P. expansum*, *P. radicicola*, *P. verrucosum*). OTA was frequently found in foodstuffs of both plant origin (e.g., cereal products, coffee, vegetable, liquorice, raisins, wine) and animal origin (e.g., pork/poultry). CIT was also found in foodstuffs of vegetable origin (e.g., cereals, pomaceous fruits, black olive, roasted nuts, spices), food supplements based on rice fermented with red microfungi *Monascus purpureus* and in foodstuffs of animal origin (e.g., cheese).



What about reversing the task ?

Creating controlled synthetic input text
from the expected relations

Aspergillus ochraceus - Ochratoxin A
Aspergillus westerdijkiae - Ochratoxin A
Aspergillus steynii - Ochratoxin A
...
Monascus purpureus - Citrinin
Penicillium expansum - Citrinin

What is inside a PubMed entry ?

What do we need ?

> [J Nat Prod.](#) 2002 Aug;65(8):1173-6. doi: 10.1021/np020009+.

5-hydroxytryptamine-derived alkaloids from two marine sponges of the genus Hyrtios

Mostafa Salmoun¹, Christine Devijver, Désiré Daloze, Jean-Claude Braekman, Rob W M van Soest

Affiliations + expand

PMID: 12193025 DOI: [10.1021/np020009+](#)

Abstract

Indonesian specimens of the marine sponges *Hyrtios erectus* and *H. reticulatus* were found to contain 5-hydroxytryptamine-derived alkaloids. Their structures were determined on the basis of their spectral properties. *H. erectus* contained hyrtiosulawesine (4), a new beta-carboline alkaloid, together with the already known alkaloids 5-hydroxyindole-3-carbaldehyde (1), hyrtiosin B (2), and 5-hydroxy-3-(2-hydroxyethyl)indole (3). *H. reticulatus* contained the novel derivative 1,6-dihydroxy-1,2,3,4-tetrahydro-beta-carboline (11) together with serotonin (5), 6-hydroxy-1-methyl-1,2,3,4-tetrahydro-beta-carboline (7), and 6-hydroxy-3,4-dihydro-1-oxo-beta-carboline (9).

MeSH terms

- > [Animals](#)
- > [Chromatography, Thin Layer](#)
- > [Indole Alkaloids / chemistry](#)
- > [Indole Alkaloids / isolation & purification*](#)
- > [Indonesia](#)
- > [Molecular Structure](#)
- > [Nuclear Magnetic Resonance, Biomolecular](#)
- > [Porifera / chemistry*](#)
- > [Serotonin / analogs & derivatives*](#)
- > [Serotonin / chemistry](#)
- > [Serotonin / isolation & purification*](#)
- > [Spectrophotometry, Ultraviolet](#)
- > [Stereoisomerism](#)

← *A title*

← *The abstract (What we want to generate)*

← *Some keywords / keyphrases*

What is inside a PubMed entry ?

What do we need ?

> J Nat Prod. 2002 Aug;65(8):1173-6. doi: 10.1021/np020009+.

5-hydroxytryptamine-derived alkaloids from two marine sponges of the genus Hyrtios

Mostafa Salmoun¹, Christine Devijver, Désiré Daloze, Jean-Claude Braekman, Rob W M van Soest

Affiliations + expand

PMID: 12193025 DOI: 10.1021/np020009+

Abstract

Given a title, some keywords and the main findings (expected relations), create a scientific abstract

← *A title*

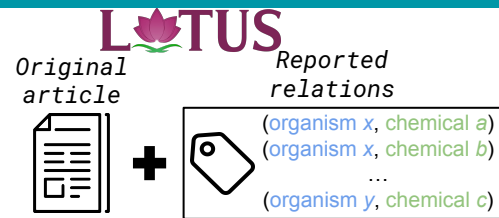
← *The abstract (What we want to generate)*

MeSH terms

- > Animals
- > Chromatography, Thin Layer
- > Indole Alkaloids / chemistry
- > Indole Alkaloids / isolation & purification*
- > Indonesia
- > Molecular Structure
- > Nuclear Magnetic Resonance, Biomolecular
- > Porifera / chemistry*
- > Serotonin / analogs & derivatives*
- > Serotonin / chemistry
- > Serotonin / isolation & purification*
- > Spectrophotometry, Ultraviolet
- > Stereoisomerism

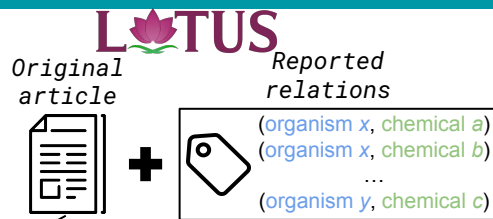
← *Some keywords / keyphrases*


Synthetic generation



Synthetic generation

```
> Extract a comma-separated list of keywords from  
the following abstract of a scientific article.  
  
----- Title -----  
-----  
----- Abstract -----  
-----  
----- kw1, kw2, ..., kwn -----  
  
output: kw1, kw2, ..., kwm 🔊  
PubChem — exclusion-list — PubTator
```



 Keyword Extraction

Not all articles have MeSH


Synthetic generation

> Extract a comma-separated list of keywords from the following abstract of a scientific article.


----- Title -----

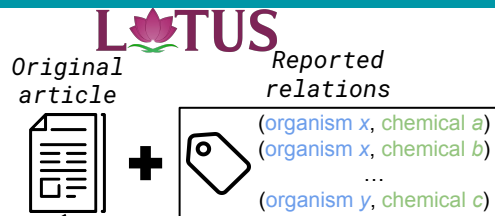
----- Abstract -----

-----kw₁, kw₂, ..., kw_n-----




output: kw₁, kw₂, ..., kw_m 

PubChem — exclusion-list — PubTator


 Keyword Extraction



(*Cystoseira usneoides*, Cystodione A)
(*Cystoseira usneoides*, Cystodione B)
...
(*Cystoseira usneoides*, Amentadione-1'-Methyl Ether)
(*Cystoseira usneoides*, Usneoidone E)

Class conversion: 
Derivatives contraction: 
Numbering: 

Main findings: Five Meroterpenoids and cystodiones A-D were isolated from *Cystoseira usneoides*.

 Findings verbaliser

Synthetic generation

```
> Extract a comma-separated list of keywords from
the following abstract of a scientific article.

----- Title -----
----- Abstract -----
----- kw1, kw2, ..., kwn -----

output: kw1, kw2, ..., kwm 📢
PubChem — exclusion-list — PubTator
```

LOTUS
Original article + Reported relations

(organism x, chemical a)
(organism x, chemical b)
...
(organism y, chemical c)

(*Cystoseira usneoides*, Cystodione A)
(*Cystoseira usneoides*, Cystodione B)
...
(*Cystoseira usneoides*, Amentadione-1'-Methyl Ether)
(*Cystoseira usneoides*, Usneoidone E)

Class conversion: ✅
Derivatives contraction: ✅
Numbering: ❌

Main findings: Five Meroterpenoids and cystodiones A-D were isolated from *Cystoseira usneoides*.

🐱 Keyword Extraction

⚙️ Instructions builder

⚙️ Findings verbaliser

```
Instructions: Given a title, a list of keywords and main findings,
create an abstract for a scientific article.
Title: Antioxidant and anti-inflammatory meroterpenoids from the
brown alga Cystoseira usneoides.
Keywords: proinflammatory cytokine tnfr, anti-inflammatory assays,
radical-scavenging activity, etc.
Main findings: Five Meroterpenoids and cystodiones A-D were isolated
from Cystoseira usneoides.
Abstract: [🐱 completing ...]
```

Synthetic generation

```
> Extract a comma-separated list of keywords from
the following abstract of a scientific article.

----- Title -----
----- Abstract -----
----- kw1, kw2, ..., kwn -----

output: kw1, kw2, ..., kwm
PubChem — exclusion-list — PubTator
```

LOTUS
Original article + Reported relations

```
(organism x, chemical a)
(organism x, chemical b)
...
(organism y, chemical c)
```

```
(Cystoseira usneoides, Cystodione A)
(Cystoseira usneoides, Cystodione B)
...
(Cystoseira usneoides, Amentadione-1'-Methyl Ether)
(Cystoseira usneoides, Usneoidone E)
```

Class conversion: ✓
Derivatives contraction: ✓
Numbering: ✗

Main findings: Five Meroterpenoids and cystodiones A-D were isolated from Cystoseira usneoides.

Keyword Extraction

Instructions builder

Findings verbaliser

Abstract Generator

```
Instructions: Given a title, a list of keywords and main findings,
create an abstract for a scientific article.
Title: Antioxidant and anti-inflammatory meroterpenoids from the
brown alga Cystoseira usneoides.
Keywords: proinflammatory cytokine tnfr, anti-inflammatory assays,
radical-scavenging activity, etc.
Main findings: Five Meroterpenoids and cystodiones A-D were isolated
from Cystoseira usneoides.
Abstract: [🐱 completing ...]
```

Synthetic generation

```
> Extract a comma-separated list of keywords from the following abstract of a scientific article.
```

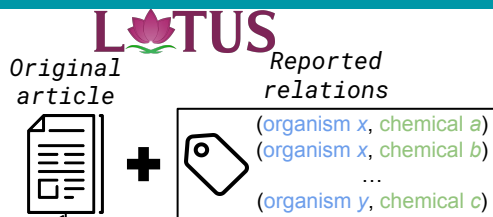
----- Title -----

----- Abstract -----

----- kw₁, kw₂, ..., kw_n -----

output: kw₁, kw₂, ..., kw_m

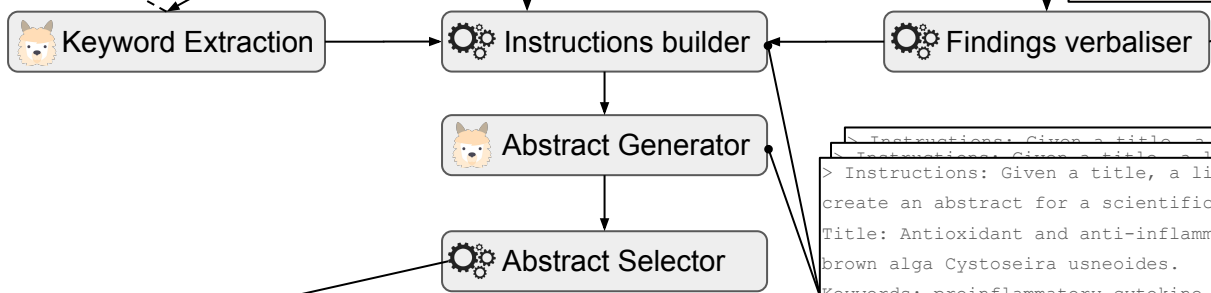
PubChem — exclusion-list — PubTator



(*Cystoseira usneoides*, Cystodione A)
(*Cystoseira usneoides*, Cystodione B)
...
(*Cystoseira usneoides*, Amentadione-1'-Methyl Ether)
(*Cystoseira usneoides*, Usneoidone E)

Class conversion: ✓
Derivatives contraction: ✓
Numbering: ✗

Main findings: Five Meroterpenoids and cystodiones A-D were isolated from *Cystoseira usneoides*.



```
Instructions: Given a title, a list of keywords and main findings, create an abstract for a scientific article.
Title: Antioxidant and anti-inflammatory meroterpenoids from the brown alga Cystoseira usneoides.
Keywords: proinflammatory cytokine tnfr, anti-inflammatory assays, radical-scavenging activity, etc.
Main findings: Five Meroterpenoids and cystodiones A-D were isolated from Cystoseira usneoides.
Abstract: [🐱 completing ...]
```

Abstract: The brown alga *Cystoseira usneoides* has been chemically studied to isolate and identify the bioactive compounds present in its tissues. In this study, five meroterpenoids and cystodiones A-D were isolated from the alga. The structures of these compounds were elucidated using spectroscopic techniques ...

- Output labels: (*Cystoseira usneoides*, Meroterpenoids); (*Cystoseira usneoides*, Cystodione A); (*Cystoseira usneoides*, Cystodione B); (*Cystoseira usneoides*, Cystodione C); (*Cystoseira usneoides*, Cystodione D)

- q = 1

Which one is synthetic ?

> J Nat Prod. 2009 Jul;72(7):1361-3. doi: 10.1021/np900181q.

Salvinorins J from *Salvia divinorum*: mutarotation in the neoclerodane system

Lukasz M Kutrzeba¹, Daneel Ferreira, Jordan K Zjawiony

Affiliations + expand

Abstract

A search for biosynthetic precursors of salvinorin A (1) led to the isolation of a new neoclerodane diterpenoid hemiacetal mixture, salvinorins J (2), from the chloroform extract of *Salvia divinorum*. A leaf surface extraction method was used on *S. divinorum*, affording a chlorophyll-free extract containing predominantly neoclerodane diterpenoids, including the new salvinorins J (2) and 14 known analogues. Salvinorins J (2) represent an example of a neoclerodane hemiacetal (lactol) susceptible to mutarotation with the formation of an equilibrium mixture of C-17 epimers.

> J Nat Prod. 2009 Jul;72(7):1361-3. doi: 10.1021/np900181q.

Salvinorins J from *Salvia divinorum*: mutarotation in the neoclerodane system

Lukasz M Kutrzeba¹, Daneel Ferreira, Jordan K Zjawiony

Affiliations + expand

Abstract

The neoclerodane system of *Salvia divinorum* has been found to undergo mutarotation, resulting in the formation of c-17 epimers. This was confirmed through the isolation and identification of several compounds from the plant using a leaf surface extraction method. The compounds isolated include Salvinorin A (1), Salvinorin B (2), Salvinorins F-H (3-5), (-)-Hardwickiic acid (6) and divinatorins A-C (7-9). These findings provide new insights into the chemical diversity of *Salvia divinorum* and highlight the importance of considering mutarotation in the study of neoclerodane systems.

Which one is synthetic ?



Real

> J Nat Prod. 2009 Jul;72(7):1361-3. doi: 10.1021/np900181q.

Salvinorins J from *Salvia divinorum*: mutarotation in the neoclerodane system

Lukasz M Kutrzeba¹, Daneel Ferreira, Jordan K Zjawiony

Affiliations + expand

Abstract

A search for biosynthetic precursors of salvinorin A (1) led to the isolation of a new neoclerodane diterpenoid hemiacetal mixture, salvinorins J (2), from the chloroform extract of *Salvia divinorum*. A leaf surface extraction method was used on *S. divinorum*, affording a chlorophyll-free extract containing predominantly neoclerodane diterpenoids, including the new salvinorins J (2) and 14 known analogues. Salvinorins J (2) represent an example of a neoclerodane hemiacetal (lactol) susceptible to mutarotation with the formation of an equilibrium mixture of C-17 epimers.



synthetic

> J Nat Prod. 2009 Jul;72(7):1361-3. doi: 10.1021/np900181q.

Salvinorins J from *Salvia divinorum*: mutarotation in the neoclerodane system

Lukasz M Kutrzeba¹, Daneel Ferreira, Jordan K Zjawiony

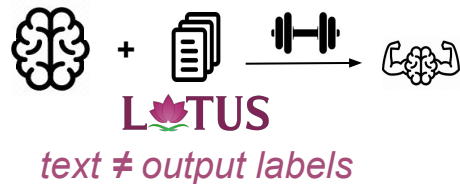
Affiliations + expand

Abstract

The neoclerodane system of *Salvia divinorum* has been found to undergo mutarotation, resulting in the formation of c-17 epimers. This was confirmed through the isolation and identification of several compounds from the plant using a leaf surface extraction method. The compounds isolated include Salvinorin A (1), Salvinorin B (2), Salvinorins F-H (3-5), (-)-Hardwickiic acid (6) and divinatorins A-C (7-9). These findings provide new insights into the chemical diversity of *Salvia divinorum* and highlight the importance of considering mutarotation in the study of neoclerodane systems.

Re-training on synthetic data

Supervised



(Weakly) supervised

Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*	
Jared Kaplan†	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Saadhvi Agarwal	Ariol Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Matenz Litwin	Scott Gray
Benjamin Chess	Jack Clark	Christopher Berner		
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei	

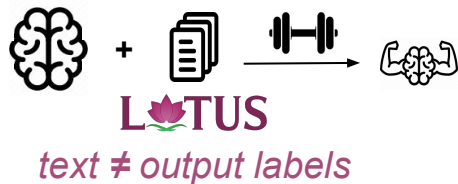
OpenAI

	precision	recall	f1
Seq2rel	47.3	5.8	10.4
GPT-2	44.8	21.7	29.3
BioGPT	42.2	26.5	32.5

LLM	precision	recall	f1
Llama-7B	27.0	9.04	13.55
Llama-13B	35.64	23.64	28.49
Llama-30B	38.51	23.24	28.99
Llama-65B	40.16	22.97	29.23
Alpaca-7B	15.14	2.21	5.86
Vicuna-13B	38.4	20.43	26.48

Re-training on synthetic data

Supervised



	precision	recall	f1
Seq2rel	47.3	5.8	10.4
GPT-2	44.8	21.7	29.3
BioGPT	42.2	26.5	32.5

(Weakly) supervised

Language Models are Few-Shot Learners			
Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*
Jared Kaplan*	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam
Amanda Askell	Saadhini Agarwal	Ariol Herbert-Voss	Gretchen Krueger
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu
Christopher Hesse	Mark Chen	Eric Sigler	Mateniz Litwin
Benjamin Chess	Jack Clark	Christopher Berner	Scott Gray
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei

OpenAI

LLM	precision	recall	f1
Llama-7B	27.0	9.04	13.55
Llama-13B	35.64	23.64	28.49
Llama-30B	38.51	23.24	28.99
Llama-65B	40.16	22.97	29.23
Alpaca-7B	15.14	2.21	5.86
Vicuna-13B	38.4	20.43	26.48

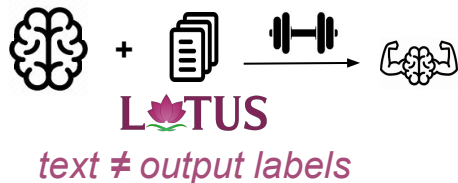
Supervised - on synthetic data



model	precision	recall	f1
seq2rel	65.1 (+17.8)	29.9 (+22.0)	41.0 (+28.9)
GPT2	52.0 (+7.2)	44.6 (+22.9)	48.0 (+18.7)
BioGPT	63.7 (+21.5)	46.5 (+20.0)	53.8 (+21.3)

Re-training on synthetic data

Supervised



	precision	recall	f1
Seq2rel	47.3	5.8	10.4
GPT-2	44.8	21.7	29.3
BioGPT	42.2	26.5	32.5

(Weakly) supervised

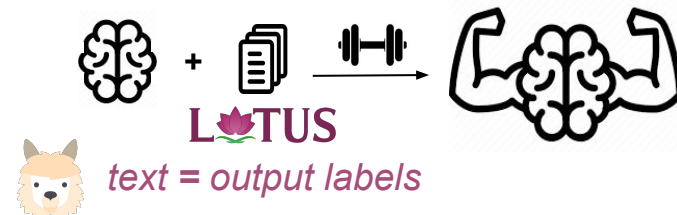
Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*
Jared Kaplan*	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam
Amanda Askell	Saadhini Agarwal	Ariol Herbert-Voss	Gretchen Krueger
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu
Christopher Hesse	Mark Chen	Eric Sigler	Mateniz Litwin
Benjamin Chess	Jack Clark	Christopher Berner	Scott Gray
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei

OpenAI

LLM	precision	recall	f1
Llama-7B	27.0	9.04	13.55
Llama-13B	35.64	23.64	28.49
Llama-30B	38.51	23.24	28.99
Llama-65B	40.16	22.97	29.23
Alpaca-7B	15.14	2.21	5.86
Vicuna-13B	38.4	20.43	26.48

Supervised - on synthetic data



model	precision	recall	f1
seq2rel	65.1 (+17.8)	29.9 (+22.0)	41.0 (+28.9)
GPT2	52.0 (+7.2)	44.6 (+22.9)	48.0 (+18.7)
BioGPT	63.7 (+21.5)	46.5 (+20.0)	53.8 (+21.3)

↓

model	precision	recall	f1
BioGPT-Large	69	51.6	59.0

Conclusion: LLM are versatile

- Fine-tuning: *garbage in, garbage out !*

Conclusion: LLM are versatile

- Fine-tuning: *garbage in, garbage out !*
- Impressive few-shot learners, but **language models** above all

Conclusion: LLM are versatile

- Fine-tuning: *garbage in, garbage out !*
- Impressive few-shot learners, but **language models** above all
- Better Synthetic Data Generator (Knowledge distillation)

Conclusion: LLM are versatile

- Fine-tuning: *garbage in, garbage out !*
- Impressive few-shot learners, but **language models** above all
- Better Synthetic Data Generator (Knowledge distillation)
- Multi-lingual opportunities

Conclusion: LLM are versatile

- Fine-tuning: *garbage in, garbage out !*
- Impressive few-shot learners, but **language models** above all
- Better Synthetic Data Generator (Knowledge distillation)
- Multi-lingual opportunities

But,

- Narrow range of “styles” compared to human-written abstracts
- Hard to control the quality
- LLM evolve fast (15 Feb. 2024 - *BioMistral*)

Thanks for your attention



André
Freitas



Magdalena
Wysocka

RELATION EXTRACTION IN UNDEREXPLORED BIOMEDICAL
DOMAINS: A DIVERSITY-OPTIMISED SAMPLING AND SYNTHETIC
DATA GENERATION APPROACH

A PREPRINT

Maxime Delmas¹, Magdalena Wysocka², and André Freitas^{1,2,3}

¹Idiap Research Institute, Switzerland

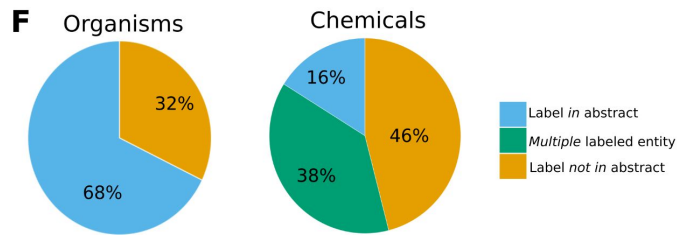
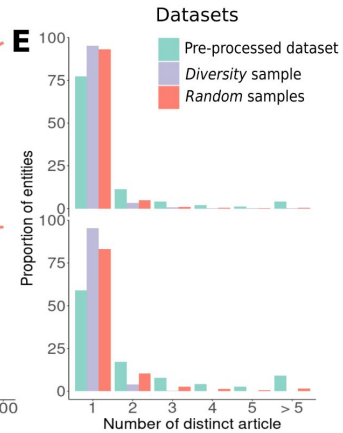
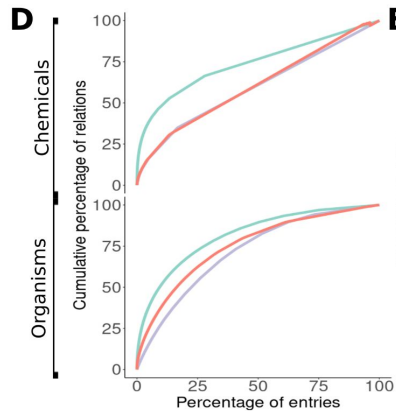
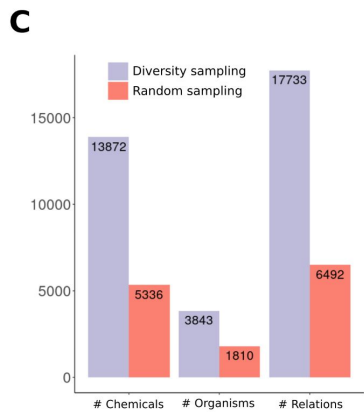
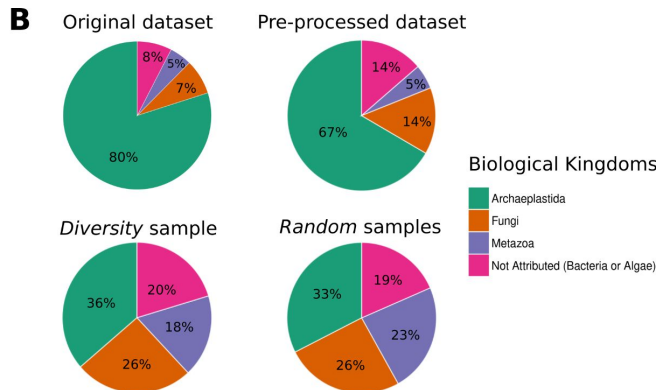
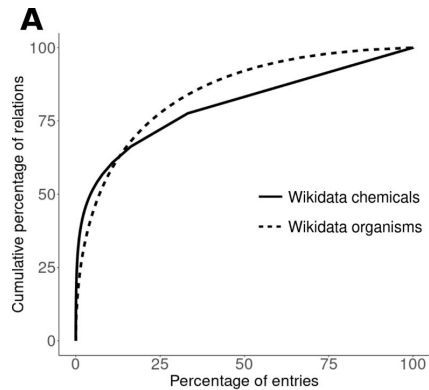
²Digital Experimental Cancer Medicine Team, Cancer Biomarker Centre, CRUK Manchester Institute

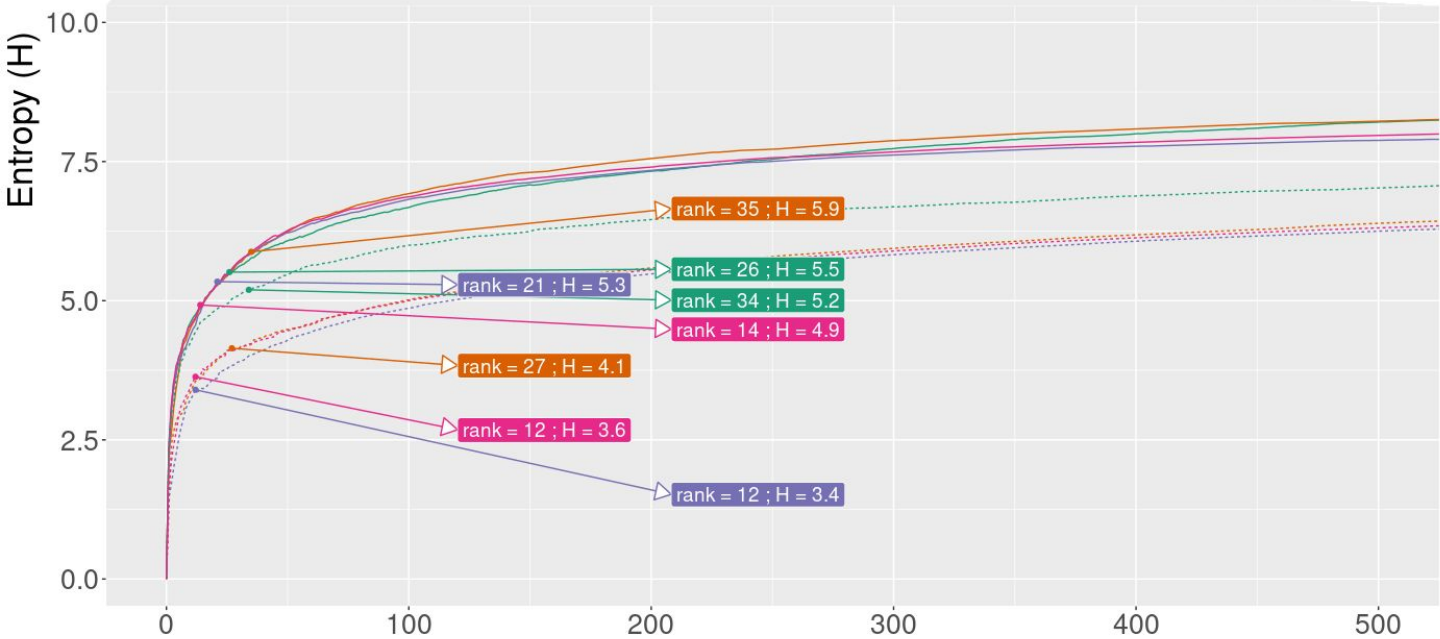
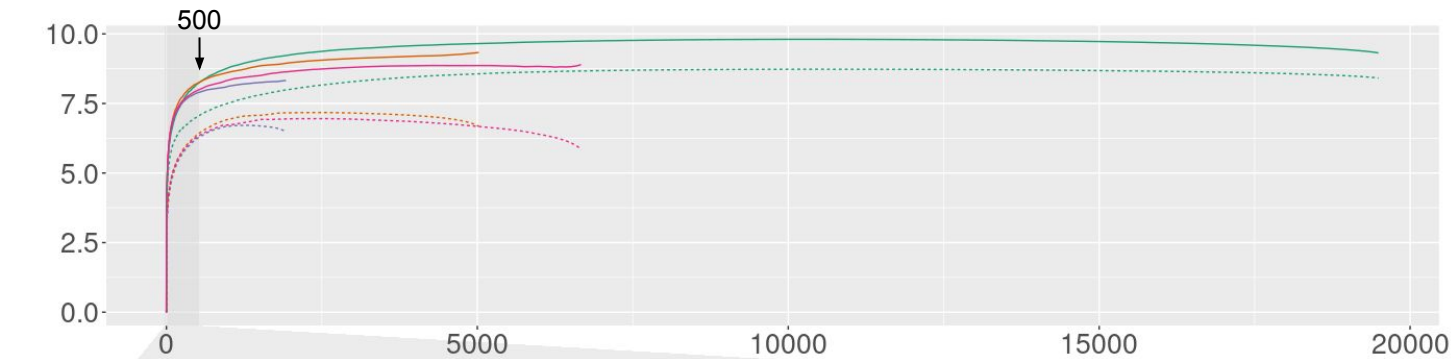
³Department of Computer Science, University of Manchester

arXiv <https://arxiv.org/pdf/2311.06364.pdf>



 <https://github.com/idiap/abroad-re>
 <https://github.com/idiap/gme-sampler>

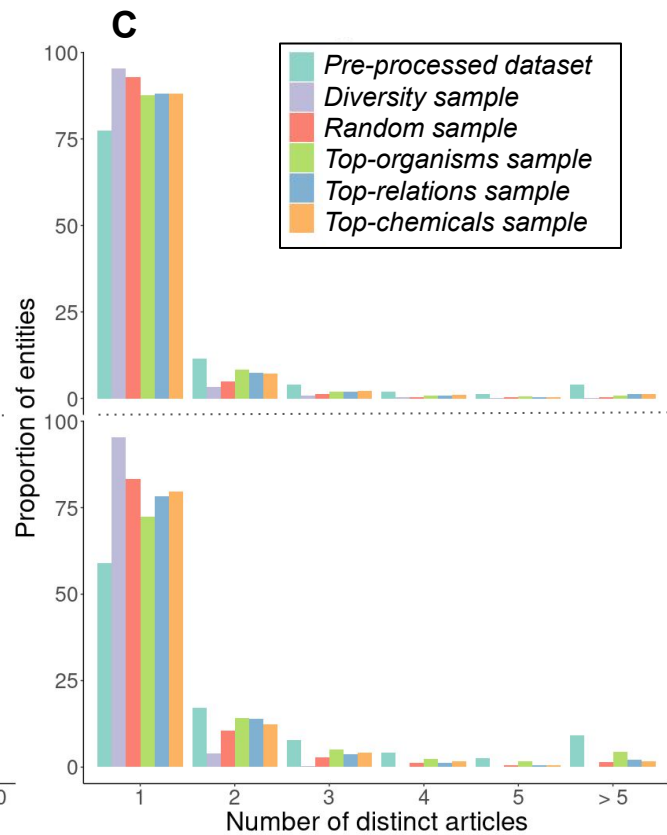
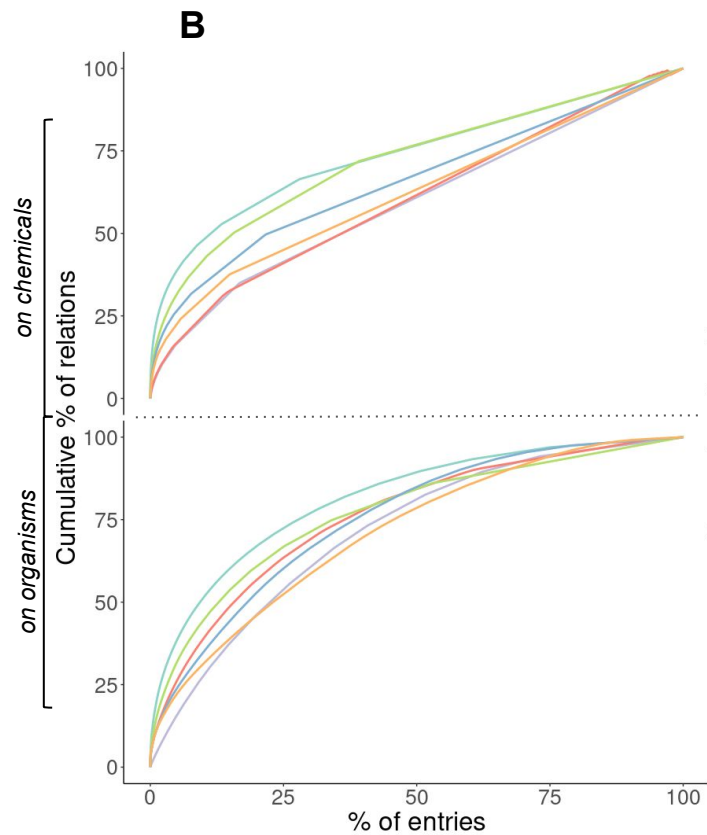
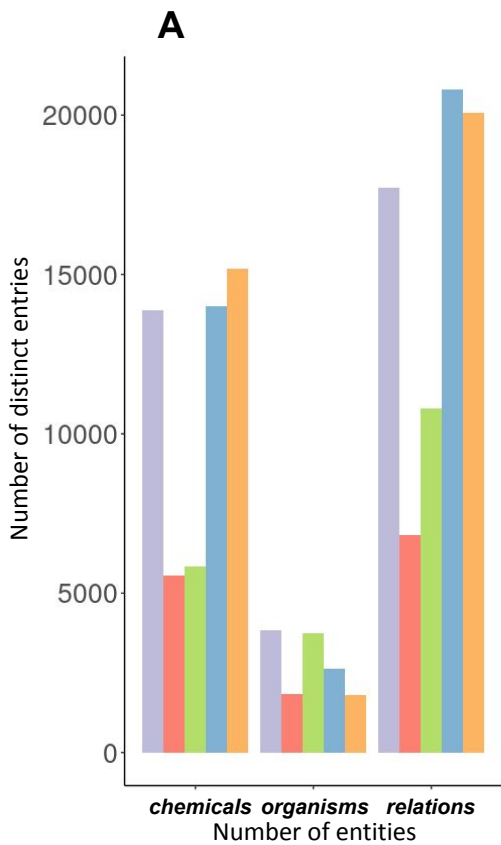


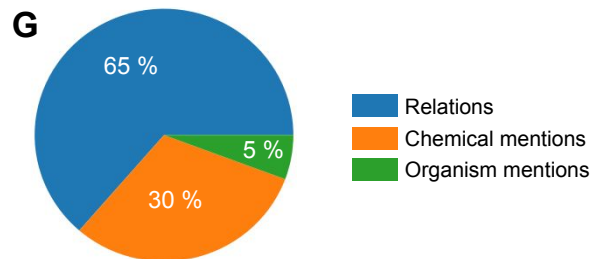
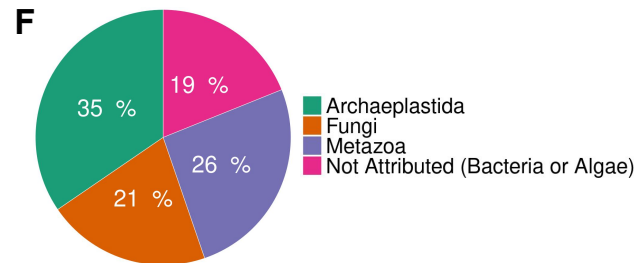
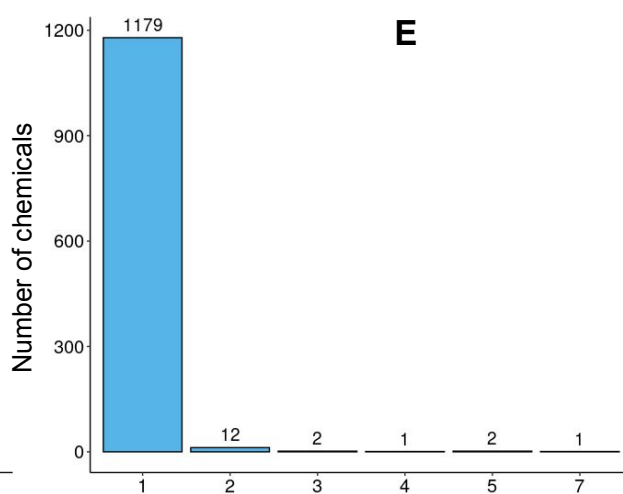
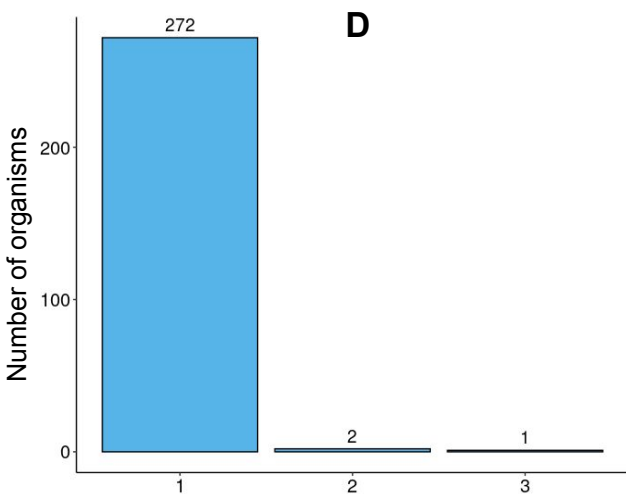
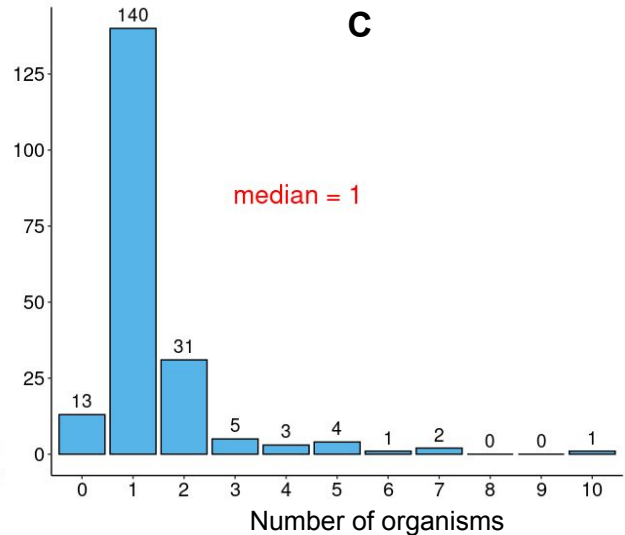
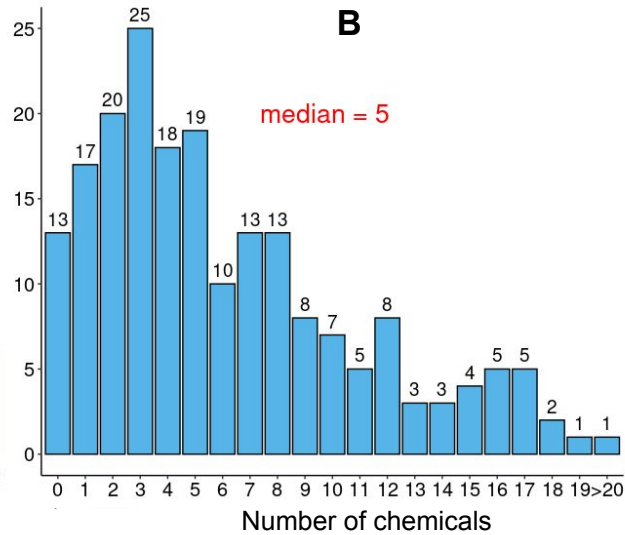
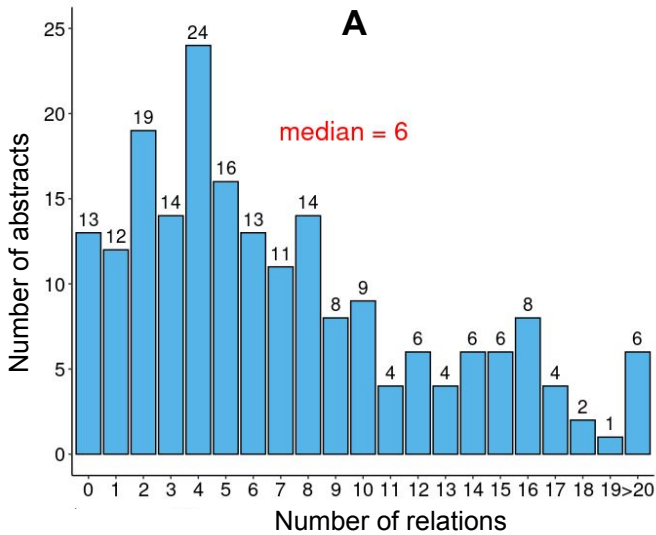


Biological Kingdoms

- Archaeplastida
- Fungi
- Metazoa
- Not Attributed (Bacteria or Algae)

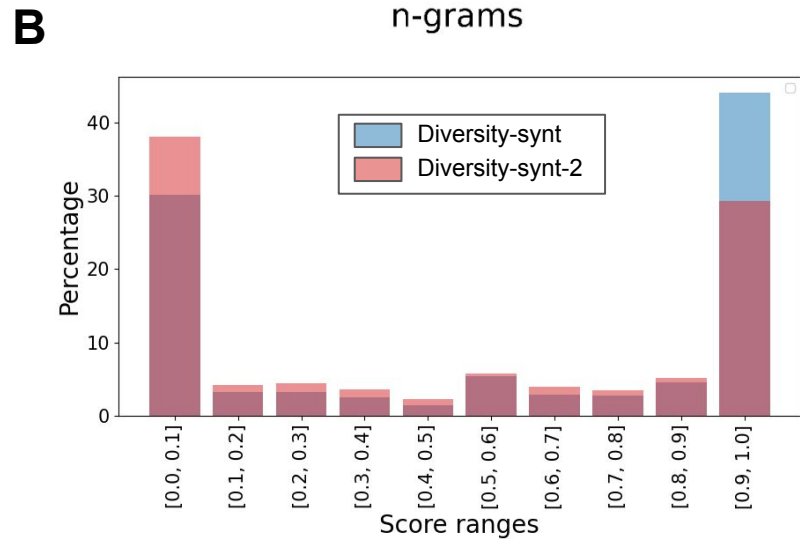
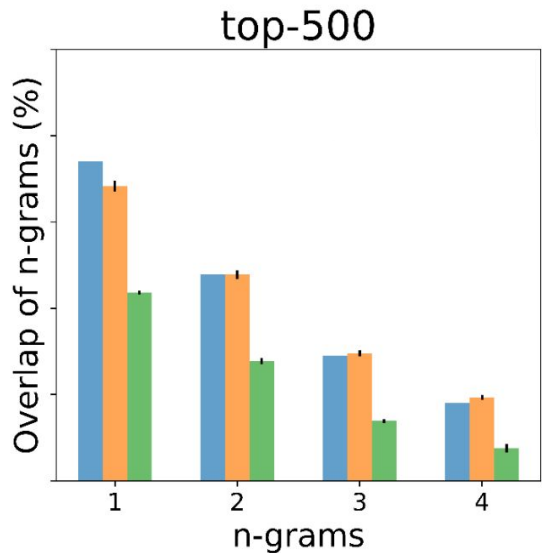
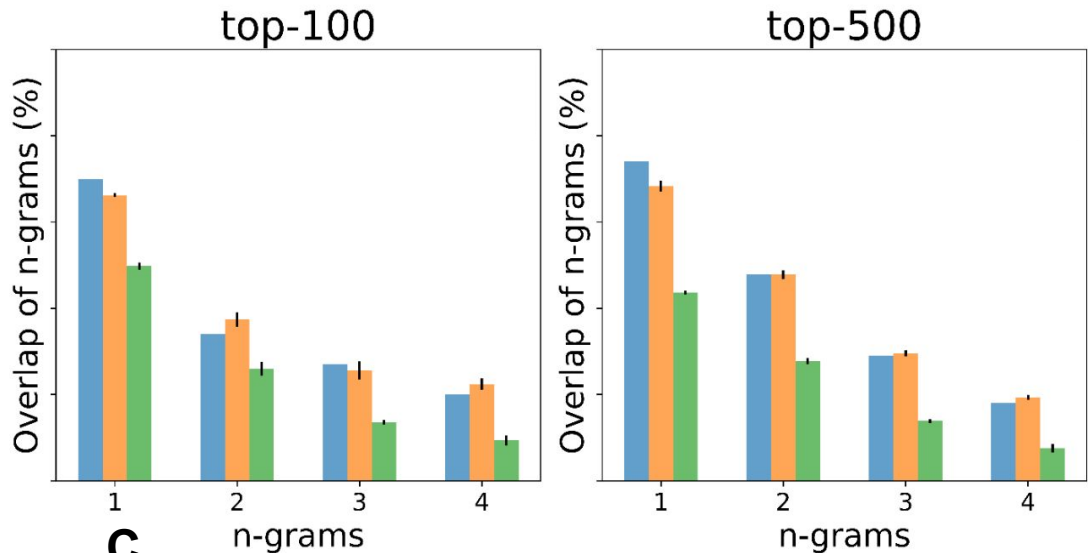
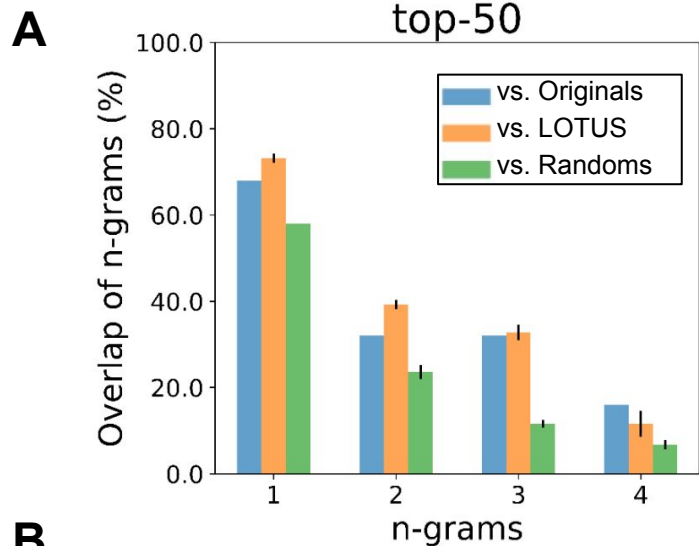
- Wikidata chemicals
- - - Wikidata organisms





model	Training	precision	recall	f1
Llama-7B	<i>Few-shots learning (5-shots)</i>	27.0	9.04	13.55
Llama-13B		35.64	23.64	28.49
Llama-30B		38.51	<u>23.24</u>	28.99
Llama-65B		<u>40.16</u>	22.97	<u>29.23</u>
Alpaca-7B		15.14	2.21	5.86
Vicuna-13B		38.4	20.43	26.48
Seq2rel	<i>Random-raw</i>	43.2 +/- (6.67)	4.8 +/- (1.16)	8.6 +/- (2.00)
	<i>Diversty-raw</i>	39.6	5.4	9.5
	<i>Extended-raw</i>	47.3	5.8	10.4
	<i>Full</i>	45.6	7.1	12.2
GPT2-QLoRA	<i>Random-raw</i>	32.5 +/- (4.83)	11.8 +/- (5.25)	15.0 +/- (2.54)
	<i>Diversty-raw</i>	22.3	19.2	20.6
	<i>Extended-raw</i>	44.8	21.7	29.3
	<i>Full</i>	47.5	22.5	30.5
BioGPT-QLoRa	<i>Random-raw</i>	47.2 +/- (4.01)	19.8 +/- (2.71)	27.6 +/- (2.48)
	<i>Diversty-raw</i>	37.1	28.4	32.2
	<i>Extended-raw</i>	42.2	26.5	32.5
	<i>Full</i>	46.7	21.3	29.3

model	Dataset	precision	recall	f1
Seq2rel	<i>Random-synt.</i>	62.4 +/- (1.03)	26.8 +/- (1.96)	37.5 +/- (1.90)
	<i>Diversty-synt.</i>	61.5	30.7	40.1
	<i>Extended-synt.</i>	65.1	29.9	41.0
GPT2-QLoRA	<i>Random-synt.</i>	42.6 +/- (2.89)	32.7 +/- (2.81)	37.2 +/- (2.80)
	<i>Diversty-synt.</i>	28.5	39.4	33.0
	<i>Extended-synt.</i>	52.0	<u>44.6</u>	<u>48.0</u>
BioGPT-QLoRa	<i>Random-synt.</i>	56.4 +/- (2.26)	38.8 +/- (1.92)	46.0 +/- 1.08
	<i>Diversty-synt.</i>	53.1	41.6	46.6
	<i>Extended-synt.</i>	<u>63.7</u>	46.5	53.8



Model	Dataset	Precision	Recall	f1
Seq2rel	<i>Diversity-synt</i>	61.5	<u>30.7</u>	<u>40.1</u>
	<i>Diversity-synt-2-selector</i>	<u>65.2</u>	28.7	39.9
	<i>Diversity-synt-2-NO-selector</i>	60.3	24.0	34.3
GPT-2	<i>Diversity-synt</i>	<u>28.5</u>	<u>39.4</u>	<u>33.0</u>
	<i>Diversity-synt-2-selector</i>	15.3	22.7	18.3
	<i>Diversity-synt-2-NO-selector</i>	15.4	20.1	17.4
BioGPT	<i>Diversity-synt</i>	<u>52.5</u>	<u>41.2</u>	<u>46.2</u>
	<i>Diversity-synt-2-selector</i>	44.5	36.4	40.0
	<i>Diversity-synt-2-NO-selector</i>	45.8	34.3	39.2