# EPFL
# RMLD

# AI generated ontologies and their impact on the quality of real-world data

Alexander Büsser, Exploris Health**
Matteo Togninalli, Isomorphic Labs*

February 12th, 2025

* Work conducted while at Visium
** Work conducted while at Idorisa

*Is Dridorexant the only sleep medication to improve daytime functioning?*

# The Challenge of Clinical Trials

- **Clinical trials are the gold standard** for evaluating safety and efficacy.

- **Challenges:**

  - Expensive and time-consuming (mean cost of cardiovascular phase 3 trials: 157M USD[1])

  - Conducted in controlled settings, limiting generalizability, operationally complex

  - Slow to adapt to real-world clinical practice needs

- **Growing interest in Real-World Evidence (RWE):**

  - Uses EHRs, insurance claims, and other clinical data

  - Can supplement or replace some aspects of traditional trials

EPFL

AMLD

*Benzodiazepine* is associated with (statistically-) *significant improvement* in *daytime functioning*

EPFL
AMLD

# Quality of Disease Representation is a key Barrier to Reproducibility

Physician Notes

Claims

Disease Representation

# Disease-Specific Medical Ontology Learning framework

EPFL
AMLD

# Disease-Specific Medical Ontology Learning framework



**Data**
- Physician notes of 82,722 insomnia patients were used for this study from Amazing Charts LLC

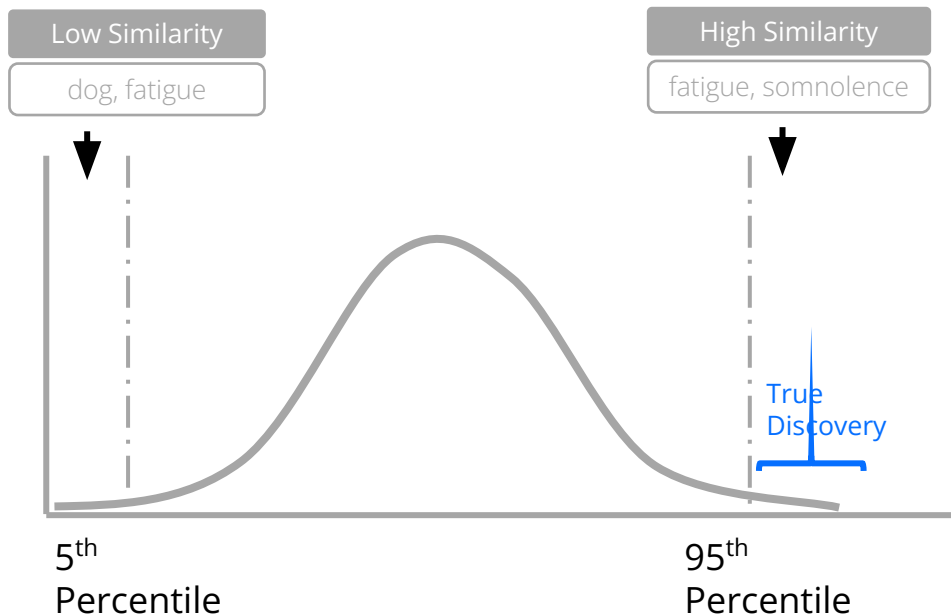**Simple ontology**
- The patient reported outcome measure Insomnia Daytime Symptoms and Impacts Questionnaire was used as the simple ontology

**Model development end evaluation**
- An ensemble of 8 word embeddings using word2vec was trained on different bootstrap samples to stabilize the concept extraction
- A evaluation metric "statistical power" was used to quantify the quality of the embedding

# A Robust Evaluation for Medical Embeddings

Statistical power measures how well embeddings distinguish true medical relationships from random ones by testing their significance against a null distribution.

**Low Similarity**

dog, fatigue

**High Similarity**

fatigue, somnolence

True Discovery

5th Percentile

95th Percentile

## How statistical power works

1. Compute **cosine similarity** for known relationships.
2. Create a **null distribution** from **random medical concept pairs**.
3. Check if known relationships score **above the 95th percentile**.
4. **Power = % of true relationships detected above chance.**
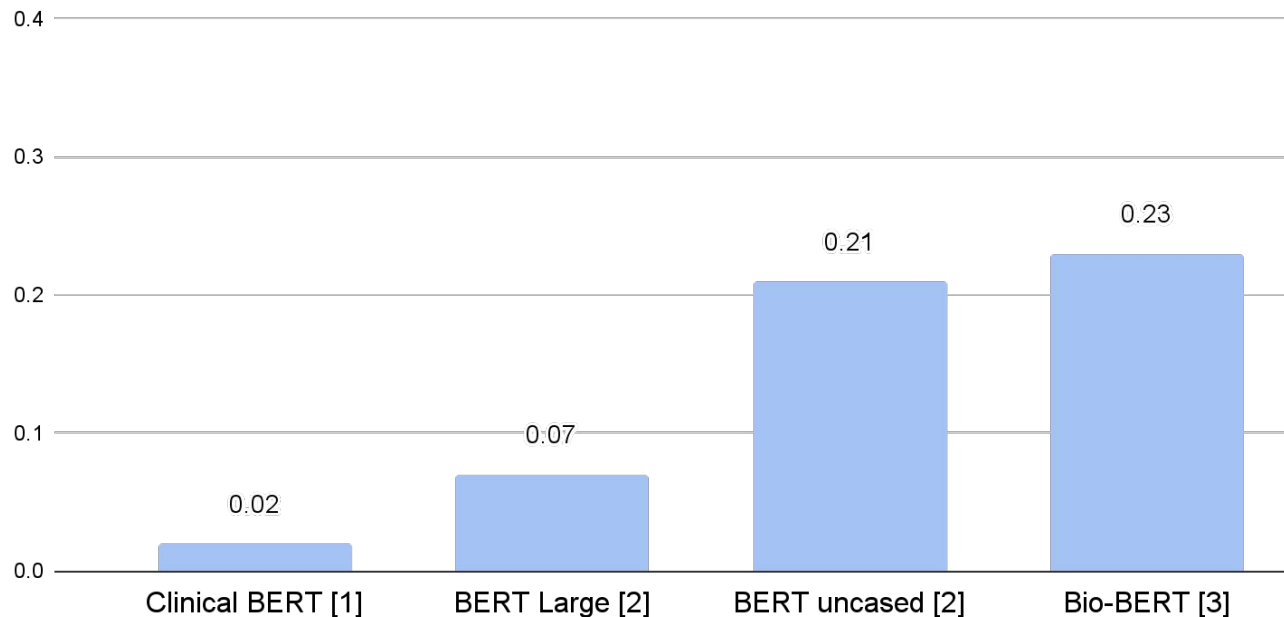
## Key advantages

- **No need for complete ground truth** – handles missing relationships.
- **Encourages discovery** – does not penalize novel relationships.

EPFL
AMLD

# Evaluation of
# word embedding

EPFL
AMLD

# Language models performance

Fraction of insomnia word-pair within **95th** percentile

[1] Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H. and Kang, J., 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. Bioinformatics, 36(4), pp.1234-1240.
[2] Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
[3] Alsentzer, E., Murphy, J.R., Boag, W., Weng, W.H., Jin, D., Naumann, T. and McDermott, M., 2019. Publicly available clinical BERT embeddings. arXiv preprint arXiv:1904.03323.

# Better than language models

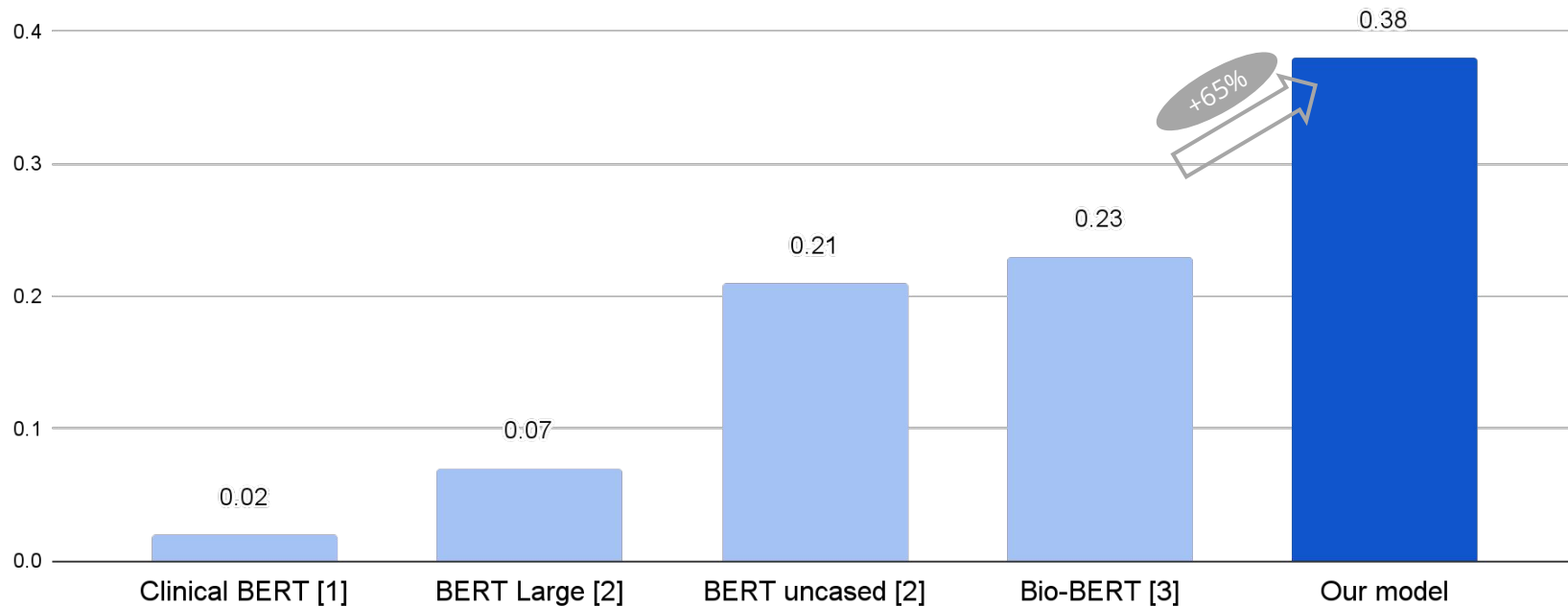Fraction of insomnia word-pair within **95th** percentile



[1] Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H. and Kang, J., 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. Bioinformatics, 36(4), pp.1234-1240.
[2] Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
[3] Alsentzer, E., Murphy, J.R., Boag, W., Weng, W.H., Jin, D., Naumann, T. and McDermott, M., 2019. Publicly available clinical BERT embeddings. arXiv preprint arXiv:1904.03323.

# Why is a simple embedding better?

- These results were obtained pre-LLM

- However, off-the-shelf LLM have struggled to obtain relevant ontologies

**Table 3.** Zero-shot results across 11 LLMs and finetuned Flan-T5-Large and Flan-T5-XL LLMs results reported for ontology learning Task A i.e. term typing in MAP@1, and as F1-score for Task B i.e. type taxonomy discovery, and Task C i.e. type non-taxonomic relation extraction. The results are in percentages.

| Task | Dataset | Zero-Shot Testing | | | | | | | | | | | Finetuned | |
| | | BERT-Large | PubMedBERT | BART-Large | Flan-T5-Large | Flan-T5-XL | BLOOM-1b7 | BLOOM-3b | GPT-3 | GPT-3.5 | LLaMA-7B | GPT-4 | Flan-T5-Large* | Flan-T5-XL* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | WordNet | 27.9 | - | 2.2 | 31.3 | 52.2 | 79.2 | 79.1 | 37.9 | **91.7** | 81.4 | 90.1 | 76.9 | **86.3** |
| | GeoNames | 38.3 | - | 23.2 | 13.2 | 33.8 | 28.5 | 28.8 | 22.4 | 35.0 | 29.5 | **43.3** | 16.9 | 18.4 |
| | NCI | 11.1 | 5.9 | 9.9 | 9.0 | 9.8 | 12.4 | 15.6 | 12.7 | 14.7 | 7.7 | **16.1** | 31.9 | **32.8** |
| | SNOMEDCT_US | 21.1 | 28.5 | 19.8 | 24.3 | 31.6 | 37.0 | **37.7** | 24.4 | 25.0 | 13.8 | 27.8 | 33.4 | **43.4** |
| | MEDCIN | 8.7 | 15.6 | 12.7 | 13.0 | 18.5 | 28.8 | **29.8** | 25.7 | 23.9 | 4.9 | 23.7 | 38.4 | **51.8** |
| B | GeoNames | 54.5 | - | 55.4 | 59.6 | 52.4 | 36.7 | 48.3 | 53.2 | **67.8** | 33.5 | 55.4 | 62.5 | 59.1 |
| | UMLS | 48.2 | 33.7 | 49.9 | 55.3 | 64.3 | 38.3 | 37.5 | 51.6 | 70.4 | 32.3 | **78.1** | 53.4 | **79.3** |
| | schema.org | 44.1 | - | 52.9 | 54.8 | 42.7 | 48.6 | 51.3 | 51.0 | **74.4** | 33.8 | 74.3 | 91.7 | **91.7** |
| C | UMLS | 40.1 | 42.7 | 42.4 | 46.0 | **49.5** | 43.1 | 42.7 | 38.8 | 37.5 | 20.3 | 41.3 | 49.1 | **53.1** |

- In the case of clinical notes, specialized terminology was not captured well by older models
- These limitations can probably be overcome by reasoning models

[3] Babaei Giglou, H., D'Souza, J., & Auer, S. (2023, October). LLMs4OL: Large language models for ontology learning. In International Semantic Web Conference (pp. 408-427)

EPFL

AMLD

# Measures of data quality

EPFL
AMLD

# Claims data with clinical experts' input

Which diagnosis codes are relevant to day-time impairment?

Search for these codes in claims data

A123
B234
C567
D123
E234
F567
G567
H567
I123
...

Clinical Expert

EPFL AMLD

# Claims with AI and clinical input

Apply AI on Physician Notes

xxxxxxxxxxxxxxxxxxxxxx**stamina**
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxx**lively**xxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxx**afraid**xxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxx**agnosia**xxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx

Share with
Clinical Experts

Which diagnosis codes are relevant to
day-time impairment?

Clinical Expert

| | |
|---|---|
| A123 | Z567 |
| B234 | D567 |
| C567 | C567 |
| D123 | W123 |
| E234 | Q567 |
| F567 | Q567 |
| G567 | O567 |
| H567 | P123 |
| I123 | ... |
| ... | ... |
| | ... |

Search for these codes in claims data

More extensive selection of codes
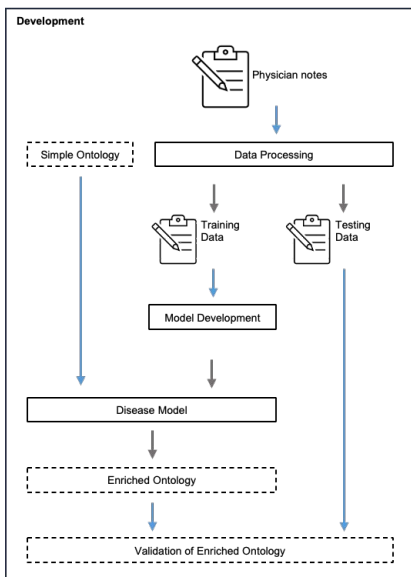
EPFL
RMLD

# Notes with AI and clinical input



AI generated ontology validated by Clinical  Expert

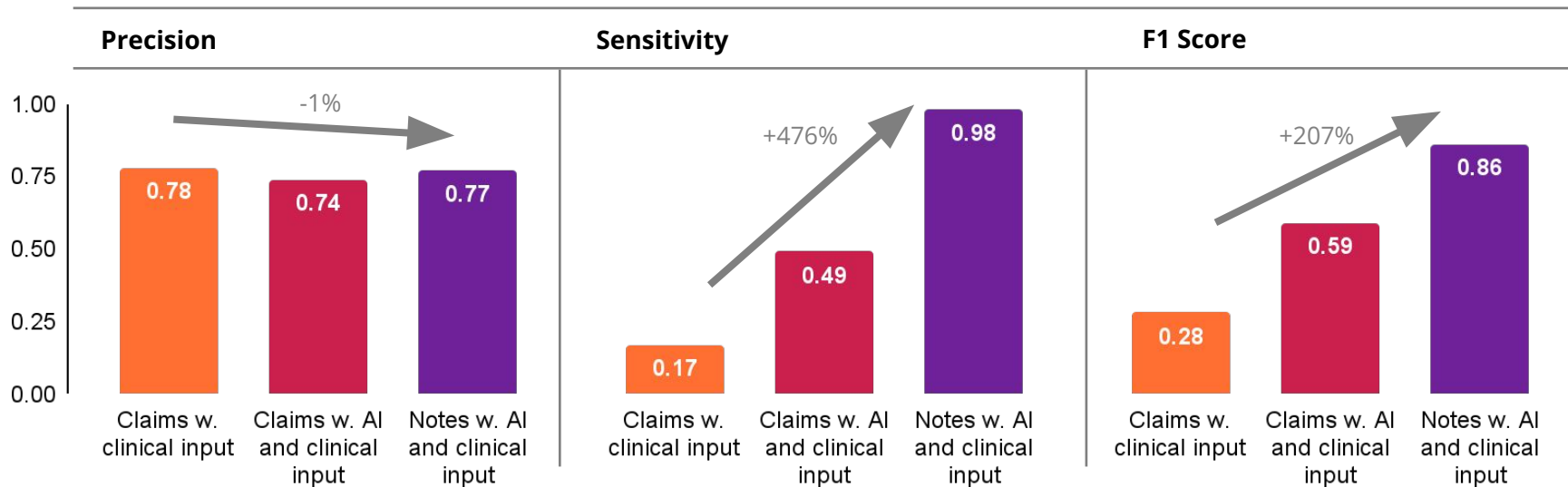Apply AI on Physician Notes

Validate with Clinical Experts

Extract from notes

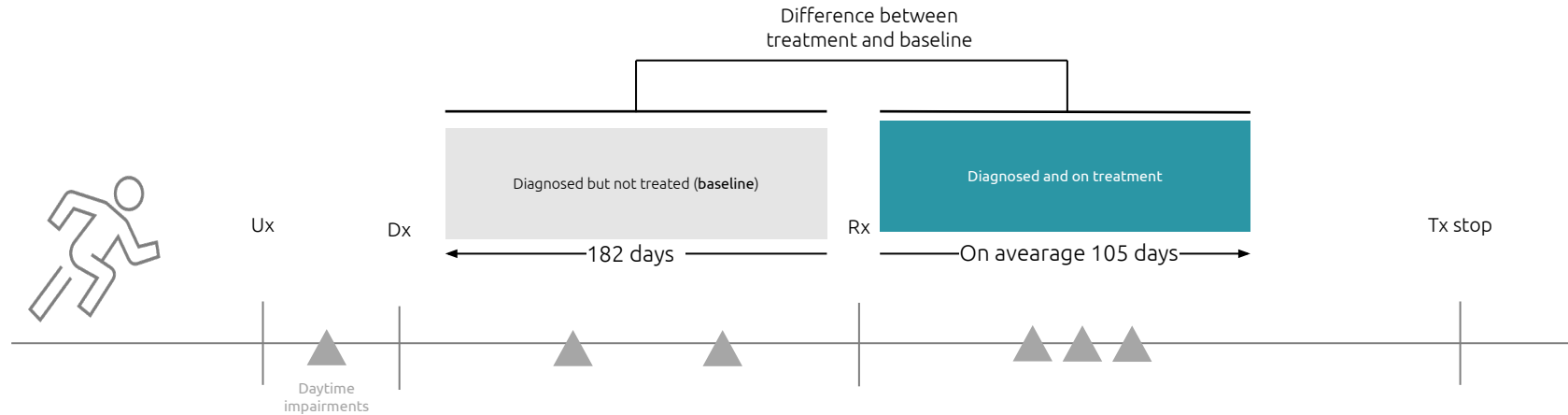Clinical Expert

# The AI-derived ontology significantly increases data quality of secondary data sources

# Real World Evidence implications

EPFL
RMLD

Study design

A within study design was chosen to control for static, non-time varying variables. The results **do not** allow for **statement of causality** between **treatments** and **daytime impairments**, only **association**

# Benzo treated patients are associated with an increased rate of day-time impairments

**Percentage increase of events per 100 patients per year** [statistically sig. numbers in **gray**] (within subject - untreated vs. Benzo treated insomnia period, n= 1045)

**Daytime impairment events ontology**

Daytime impairment
**19%**

Mood domain
**26%**

Alert/Co. domain
7%

Sleep domain
21%

# Daytime impairment events ontology

**Percentage increase of events per 100 patients per year**

[statistically sig. numbers in **gray**]

| | Treated with Trazadone | Treated with Z Drug | Treated with Benzo |
|---|---|---|---|
| **Daytime impairment** | 9 % | 16 % | 19 % |
| Mood domain | 13 % | 22 % | 26 % |
| Alert/Co. domain | 4 % | 16 % | 7 % |
| Sleep domain | 7 % | 12 % | 21 % |

EPFL
AMLD

# Conclusion

EPFL
AMLD

# Conclusion

- We have demonstrated that we can learn a representation of a disease in an **unsupervised** manner that is reflective on how a condition is describe in **real-world** setting

- By doing so we have significantly increased the data quality and thus the reliability of resulting studies

- Extracting clinically relevant endpoints form textual data remains a difficult challenge. We believe the success of this project is rooted in ...
  - setting realistic expectations on the program's outcome
  - clear focus on one indication
  - strong cross-functional collaboration between AI experts and clinical specialists,
  - rigorous validation of the system's reliability

- While traditional LLMs are not good at creating meaningful ontologies, more recent reasoning systems offer better capabilities in this space.

EPFL
AMLD

# Acknowledgments

EPFL

RMLD

# Experts involved and contributions to the Study

| | Contribution | People |
|---|---|---|
| **VISIUM** | Development of AI, data pipelines | Renato Durrer, Tobias Ochsner Moritz Freidank, Thibault Viglino Matteo Togninalli |
| **ETH** | Overall methodological validation | Prof. Karsten Borgwardt, ETH |
| **IBM** | Validation environment setup | Vlad Zamfirescu |
| Medical College of Georgia, Augusta University | Clinical validation | M.D. William Vaughn McCall |
| The University of Arizona | Clinical validation | M.D. Michael Grandner |
| **idorsia** | Clinical | Andrea Beyer |
| | Medical Affairs | Antonio Olivieri |
| | Value and Access | Paulien Meijer |
| | Communication | Andrew Jones, Agnes Lei |
| | GIS | Oliver Cotto, Thomas Straehl |

EPFL
AMLD

# Evaluation metrics

**Sensitivity**
Measures what fraction of people that have a condition are recognized as such.

**Precision**
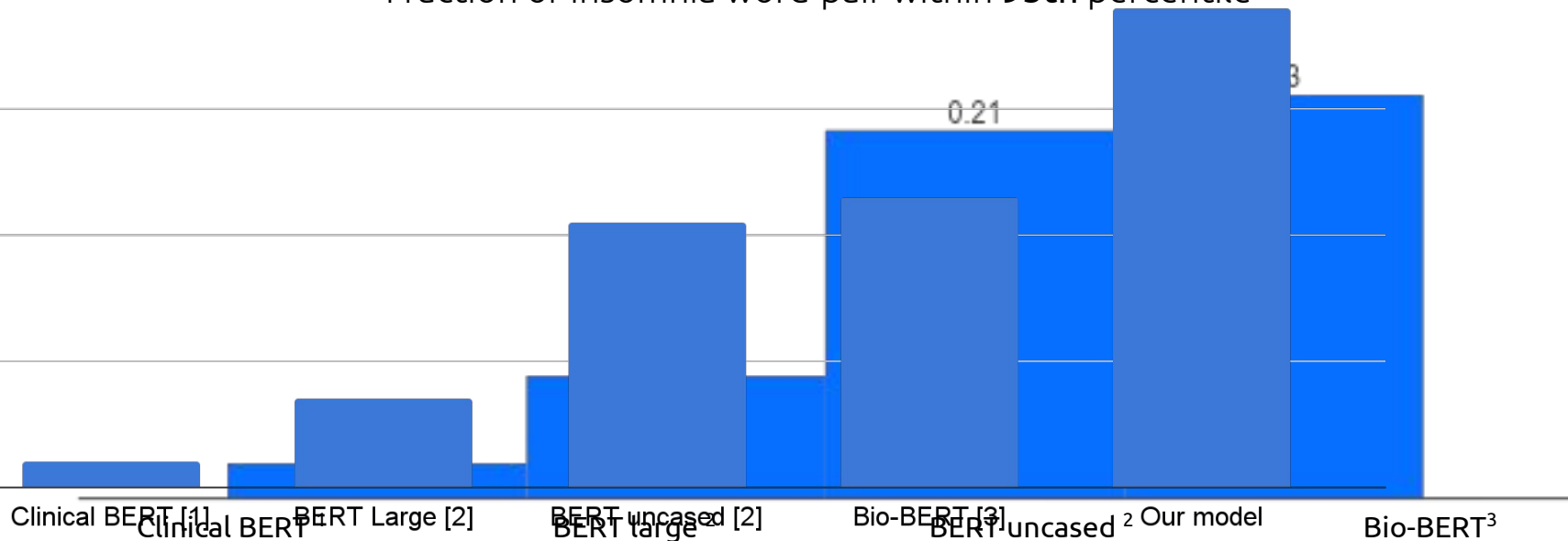Measures what fraction of people that are recognized to have a condition actually have it.

**F1 Score**
F1 is the harmonic mean between Sensitivity and precision.

A high F1 score corresponds to good **data quality**

EPFL
AMLD

# Language models performance



Fraction of insomnia word-pair within **95th** percentile

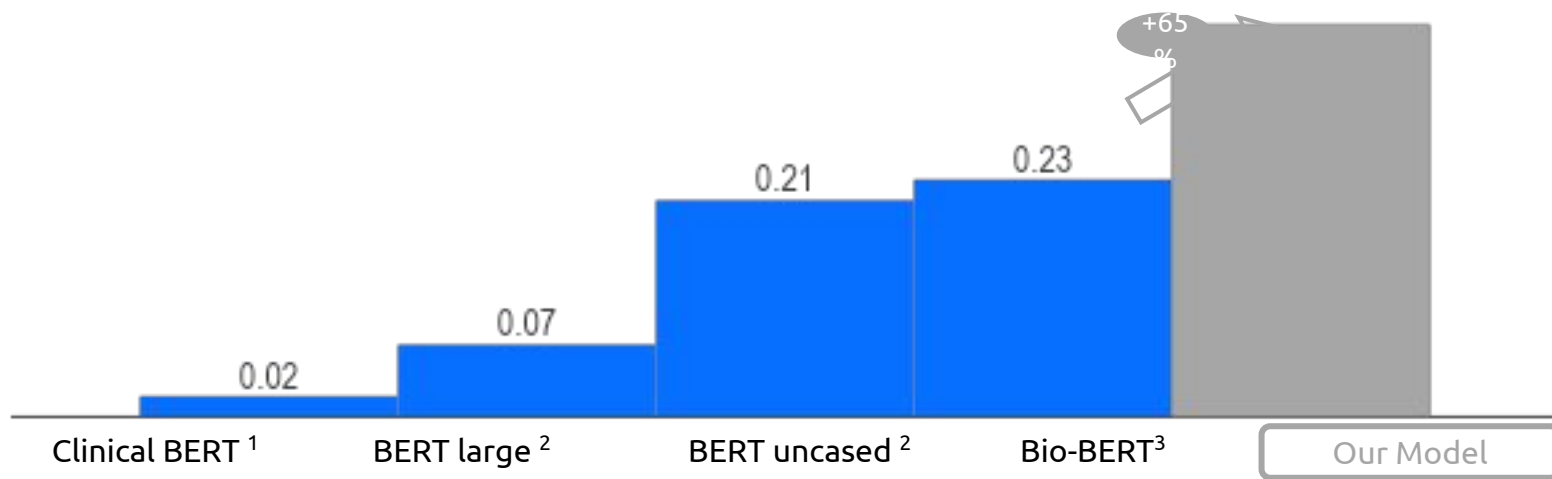Clinical BERT | BERT large | BERT uncased ² | Bio-BERT³

[1] Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H. and Kang, J., 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. Bioinformatics, 36(4), pp.1234-1240.
[2] Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
[3] Alsentzer, E., Murphy, J.R., Boag, W., Weng, W.H., Jin, D., Naumann, T. and McDermott, M., 2019. Publicly available clinical BERT embeddings. arXiv preprint arXiv:1904.03323.

EPFL
AMLD

# Better than State-of-the-art language models

## Fraction of insomnia word-pair within **95th** percentile



Clinical BERT [1] — 0.02
BERT large [2] — 0.07
BERT uncased [2] — 0.21
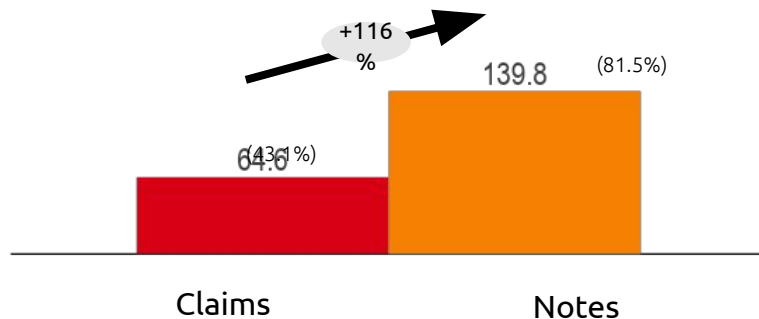Bio-BERT [3] — 0.23
Our Model — +65%

[1] Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H. and Kang, J., 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. Bioinformatics, 36(4), pp.1234-1240.
[2] Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
[3] Alsentzer, E., Murphy, J.R., Boag, W., Weng, W.H., Jin, D., Naumann, T. and McDermott, M., 2019. Publicly available clinical BERT embeddings. arXiv preprint arXiv:1904.03323.
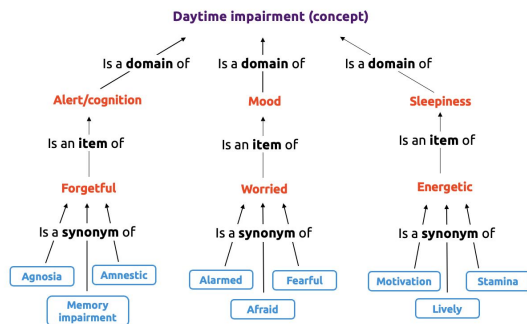
EPFL
AMLD

# The prevalence of daytime impairments has been strongly underestimated in claims data

The number of daytime impairment in **100 patient years** has been calculated for a fixed insomnia population



+116%

139.8 (81.5%)

64.6 (43.1%)

Claims

Notes

EPFL
RMLD

# Notes with AI and clinical input



AI generated ontology validated by Clinical Expert

Apply AI on Physician Notes

**Daytime impairment (concept)**

Is a **domain** of — Is a **domain** of — Is a **domain** of

Alert/cognition — Mood — Sleepiness

Is an **item** of — Is an **item** of — Is an **item** of

Forgetful — Worried — Energetic

Is a **synonym** of — Is a **synonym** of — Is a **synonym** of

Agnosia — Amnestic — Alarmed — Fearful — Motivation — Stamina

Memory impairment — Afraid — Lively

Validate with Clinical Experts

Clinical Expert

Extract from notes

stamina
lively
afraid
agnosia

EPFL
AMLD

# The AI-derived ontology significantly increase data quality of secondary data sources



**Precision**

0.78 — Claims w. clinical input
0.74 — Claims w. AI and clinical input
0.77 — Notes w. AI and clinical input
-1%

**Sensitivity**

0.17 — Claims w. clinical input
0.49 — Claims w. AI and clinical input
0.98 — Notes w. AI and clinical input
+476%

**F1 Score [1]**

0.28 — Claims w. clinical input
0.59 — Claims w. AI and clinical input
0.86 — Notes w. AI and clinical input
+207%

[1] Van Rijsbergen, C. J. (1979) Information Retrieval. London: Butterworths

EPFL
AMLD