# IDIOMATIC MWE: WHY?

**Idiomatic multiword expressions appear to be more challenging for translation (Evjen, 2018).**

WHY?

How can we detect idiomatic multiword expressions from an ontology without examples of usage in context of the phrase?

WHAT?

HOW?

BioLORD Self-Explainability Score indicates how much a semantic model believes word interaction matters to describe a phrase

5

# IDIOMATIC MWE: WHAT?

**Idiomatic multiword expressions appear to be more challenging for translation (Evjen, 2018).**

WHY?

WHAT?

HOW?

**How can we detect idiomatic multiword expressions from an ontology without examples of usage in context of the phrase?**

**BioLORD Self-Explainability Score indicates how much a semantic model believes word interaction matters to describe a phrase**

# IDIOMATIC MWE: HOW?

Idiomatic multiword expressions appear to be more challenging for translation (Evjen, 2018).

WHY?

WHAT?

HOW?

How can we detect idiomatic multiword expressions from an ontology without examples of usage in context of the phrase?

BioLORD Self-Explainability Score indicates how much a semantic model believes word interaction matters to describe a phrase

# RELATED WORK: BIOLORD

The representations of concepts and sentences produced by BioLORD are more semantic than state of the art alternatives.
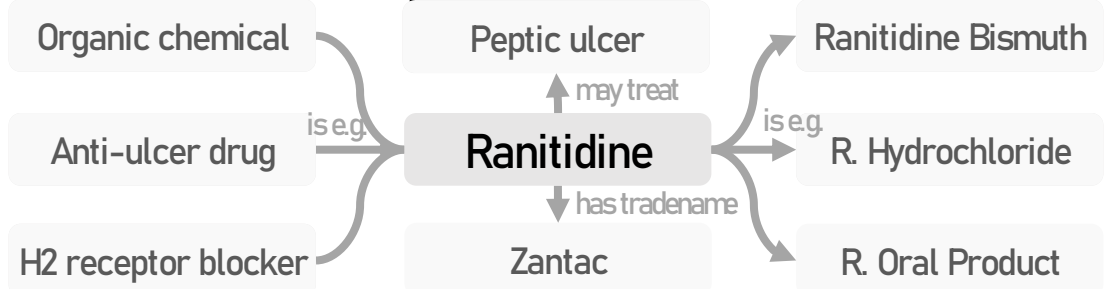
...BUT HOW?

# THE BIOLORD TRAINING

# THE BIOLORD TRAINING

**⊗ Dictionary**

**Ranitidine** ... is a non–imidazole blocker of those
= Zantac histamine receptors which can mediate
= Ranisen gastric secretion (H2 receptors). It is
= Taladine used to treat gastrointestinal ulcers.

**⊕ Knowledge Graph**

Organic chemical

Anti–ulcer drug — is e.g. →

H2 receptor blocker

Peptic ulcer

↑ may treat

**Ranitidine**

↓ has tradename

Zantac

Ranitidine Bismuth

is e.g.

R. Hydrochloride

R. Oral Product

**We propose a strategy to** →

Learn ontological
representation

**CONCEPTS**

**or**

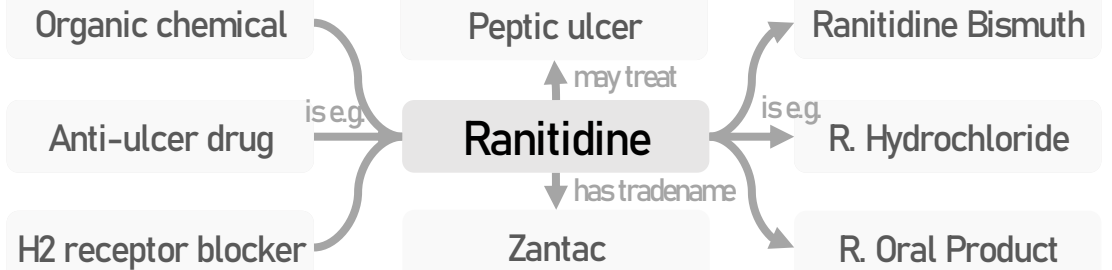**SENTENCES**

Semantic
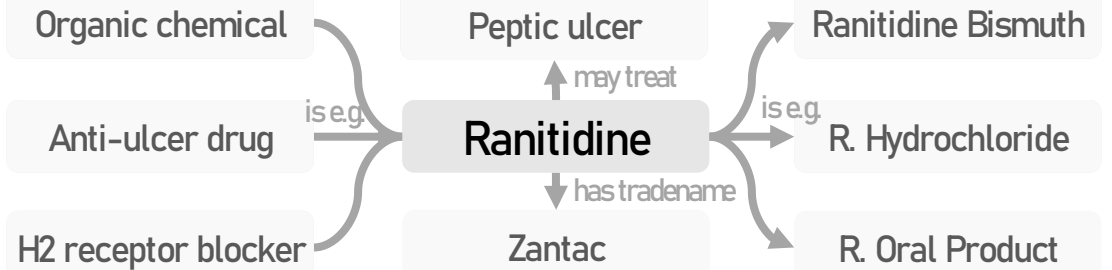model

# THE BIOLORD TRAINING

## ⊗ Dictionary

**Ranitidine**
= Zantac
= Ranisen
= Taladine

... is a non-imidazole blocker of those histamine receptors which can mediate gastric secretion (H2 receptors). It is used to treat gastrointestinal ulcers.

## ⊕ Knowledge Graph

Organic chemical — Peptic ulcer — Ranitidine Bismuth

Anti-ulcer drug — *is e.g.* — Ranitidine — *is e.g.* — R. Hydrochloride

H2 receptor blocker — *has tradename* — Zantac — R. Oral Product

*may treat*

**We propose a strategy to** →

Learn ontological representation

**CONCEPTS** / **SENTENCES**  or  → Semantic model → **EMBEDDINGS**
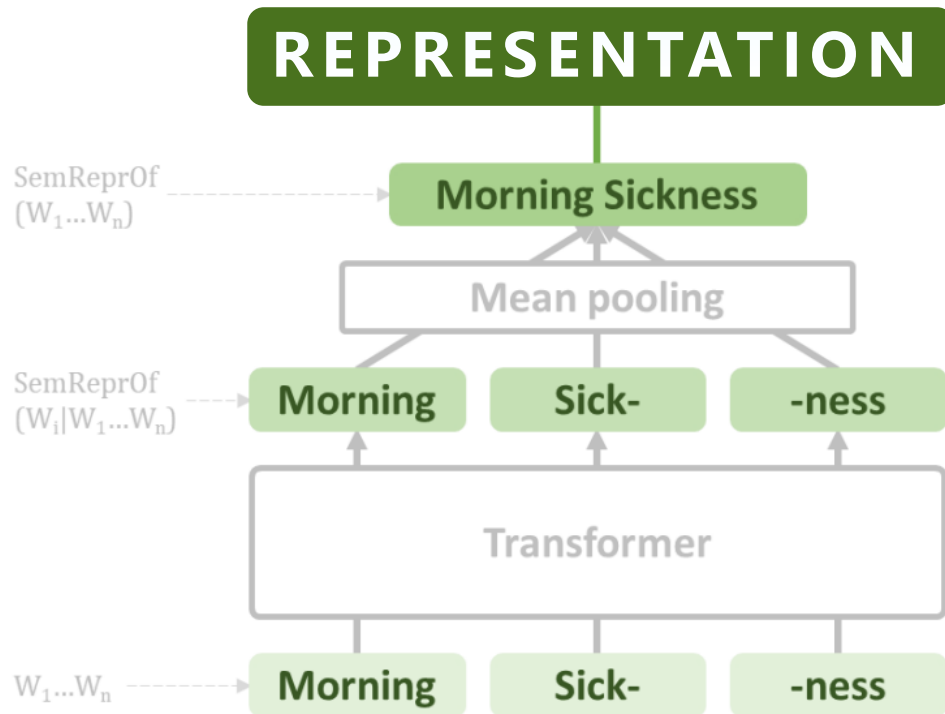
**BioLORD Self-Explainability Score indicates how much a semantic model believes word interaction matters to describe a phrase…**

**…BUT HOW?**

# SEMANTIC REPRESENTATION

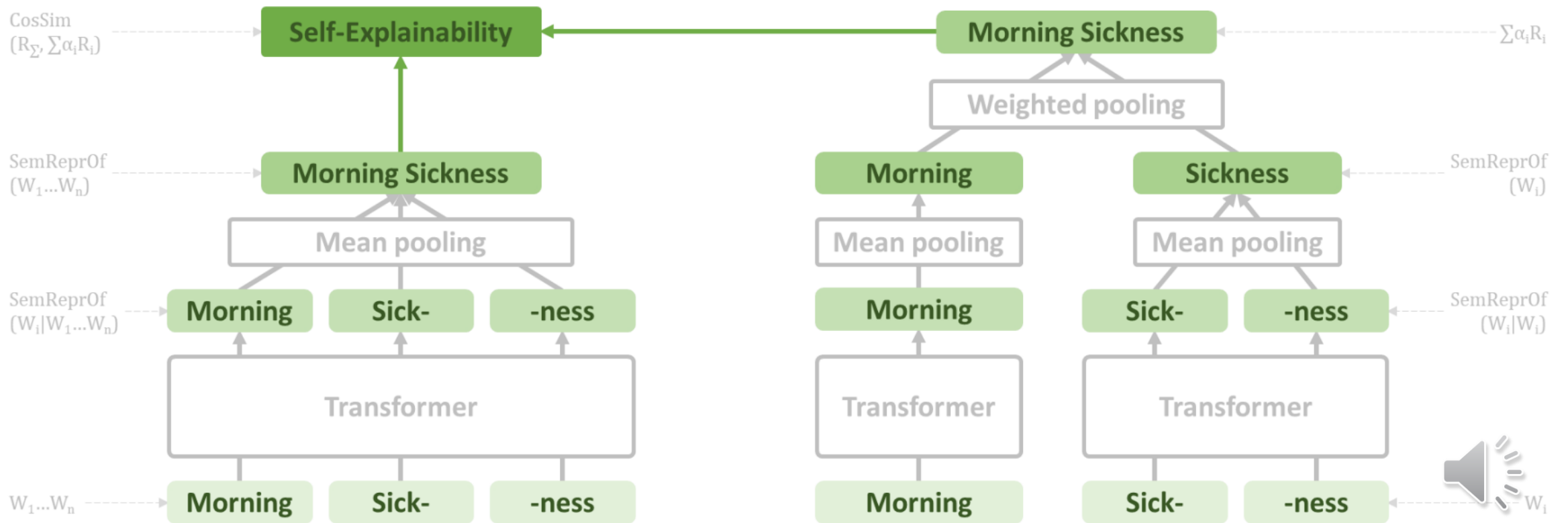**We first compute a semantic representation for the concept**



The semantic representation of concepts is computed by averaging the embedding of its tokens, after interactions between tokens have been taken into consideration by the Transformer model
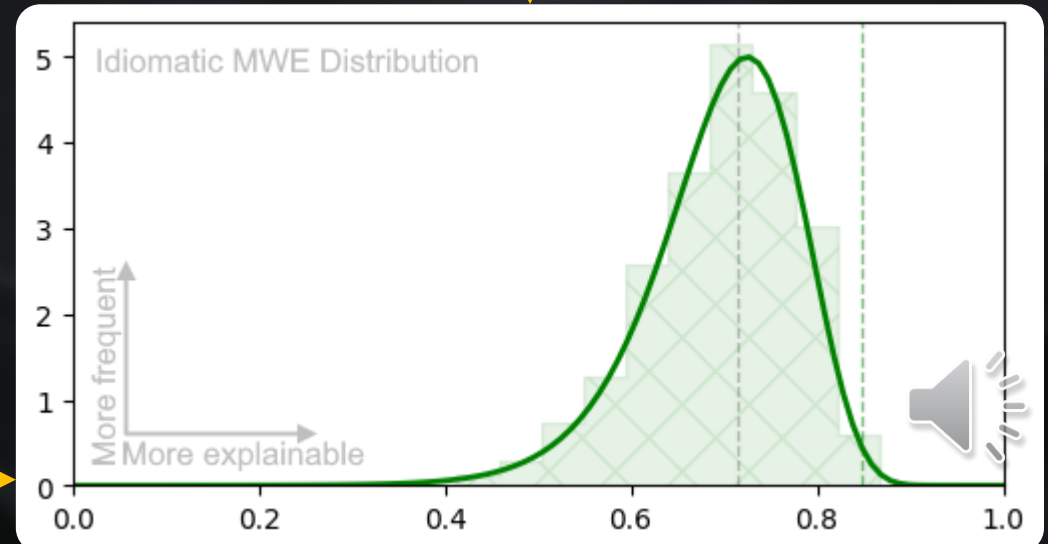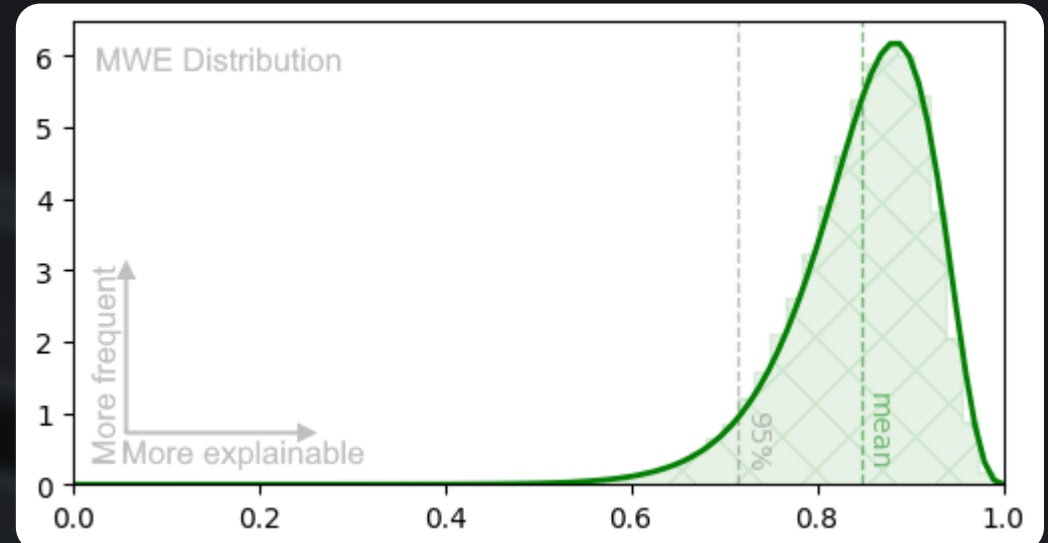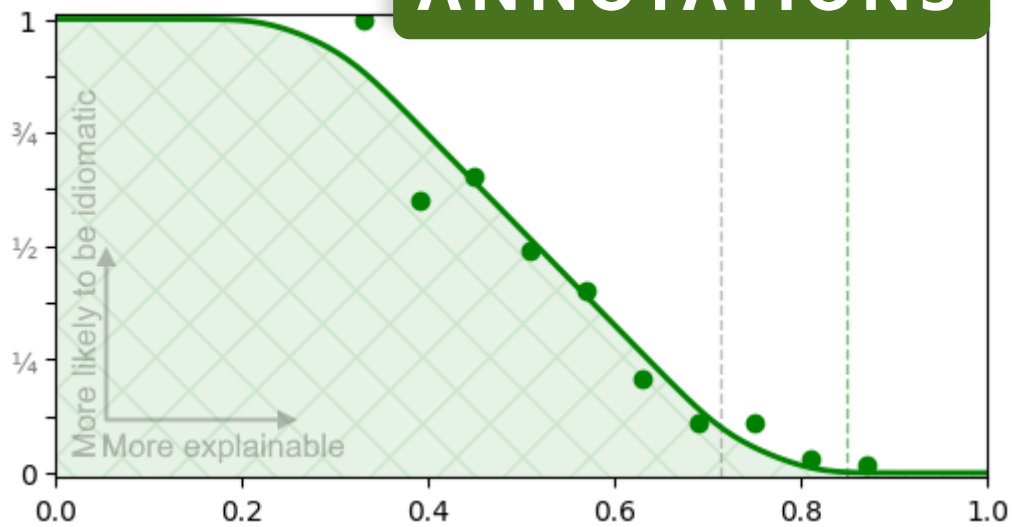
**We can compare it to an average that ignores interactions**
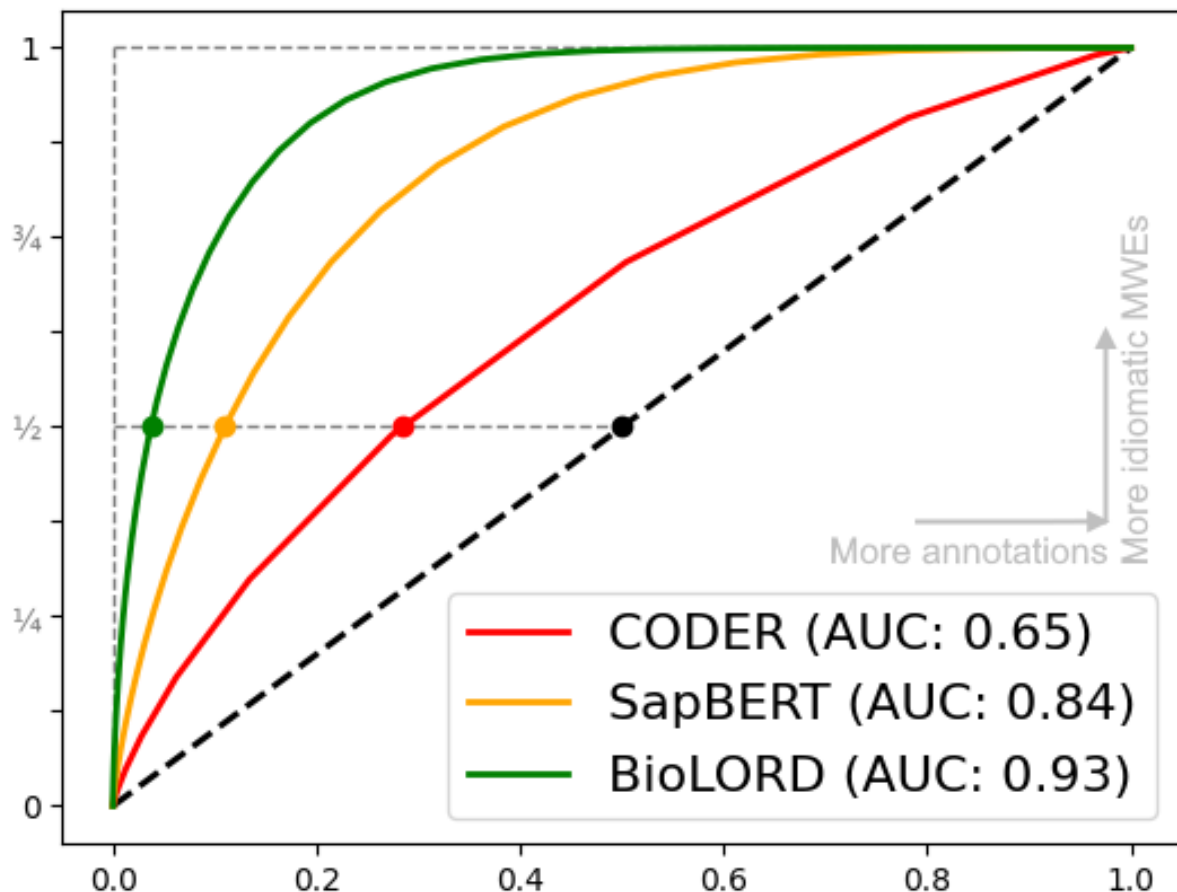
# RESULTS

**We sample the MWEs by score bucket, and classify them between self-explanatory and idiomatic...**

**ANNOTATIONS**

# RESULTS (ROC CURVE)



To recall 50% of all idiomatic expressions in the dataset, it is only required to annotate 5% of the data using the BioLORD-based score.

Other models, which do not ground their representation using definitions, perform significantly worse.

**THE END**

# Thank you

## for listening to this talk

@FremyCompany          francois.remy@ugent.be