# Building explanation machines for Science

## A Neuro-symbolic perspective
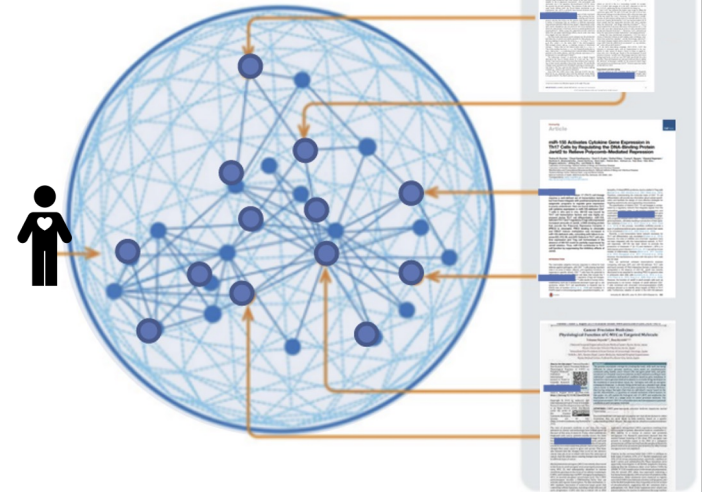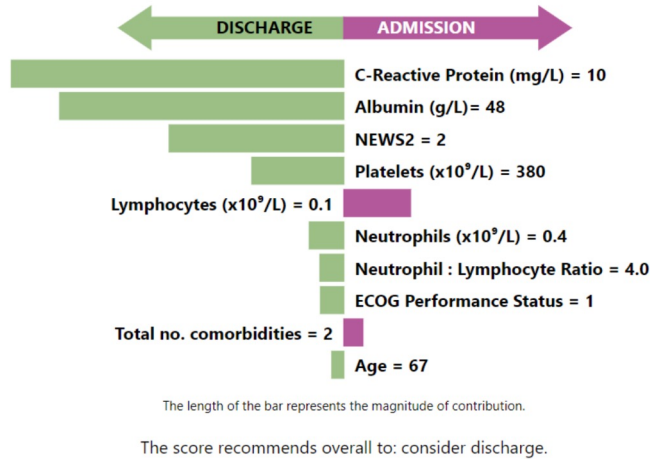
André Freitas
Reasoning & Explainable AI (ExplAIn) Lab

IPAM UCLA
(January 2023)

ExplAIn Lab

# Three Perspectives on Explanation

**Important Features Contributing to the Model Prediction for Your Patient**

DISCHARGE ← → ADMISSION

- C-Reactive Protein (mg/L) = 10
- Albumin (g/L) = 48
- NEWS2 = 2
- Platelets (x$10^9$/L) = 380
- Lymphocytes (x$10^9$/L) = 0.1
- Neutrophils (x$10^9$/L) = 0.4
- Neutrophil : Lymphocyte Ratio = 4.0
- ECOG Performance Status = 1
- Total no. comorbidities = 2
- Age = 67

The length of the bar represents the magnitude of contribution.

The score recommends overall to: consider discharge.

## Expert-AI Interaction

## Natural Language Explanations

Genes — Pathways — Biological processes

- Mutations
- Copy number
- Fusion
- Methylation
- Gene expression

G1, G2, G3, G4, G5, G6 — P1, P2, P3 — BP1, BP2, BP3 — Outcome

Patient profile → Biologically-informed architecture → Interpretation → Experimental and clinical validation

Genes: G1 G2 G3 G4 — Viability
Pathways: P1 P2 P3 — Response
Processes: BP1 BP2 BP3 — Survival

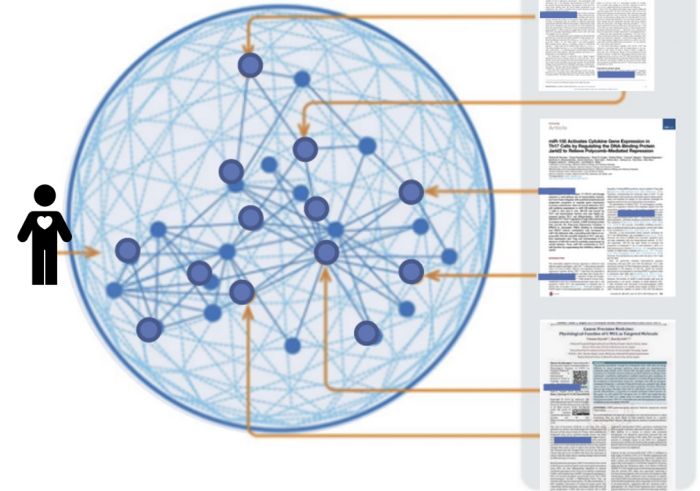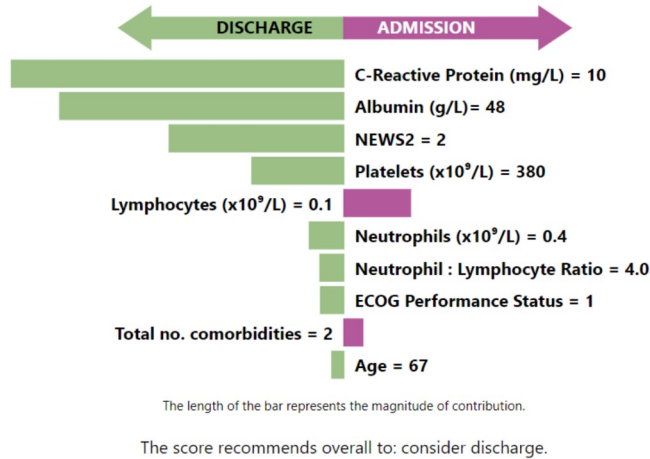## Prior knowledge & explainability

Elmarakeby et al, Nature (2021)

# Three Perspectives on Explanation

**Important Features Contributing to the Model Prediction for Your Patient**

DISCHARGE ← | → ADMISSION

C-Reactive Protein (mg/L) = 10
Albumin (g/L) = 48
NEWS2 = 2
Platelets (x10⁹/L) = 380
Lymphocytes (x10⁹/L) = 0.1
Neutrophils (x10⁹/L) = 0.4
Neutrophil : Lymphocyte Ratio = 4.0
ECOG Performance Status = 1
Total no. comorbidities = 2
Age = 67

The length of the bar represents the magnitude of contribution.

The score recommends overall to: consider discharge.

## Expert-AI Interaction

## Natural Language Explanations

## Prior knowledge & explainability

Elmarakeby et al, Nature (2021)

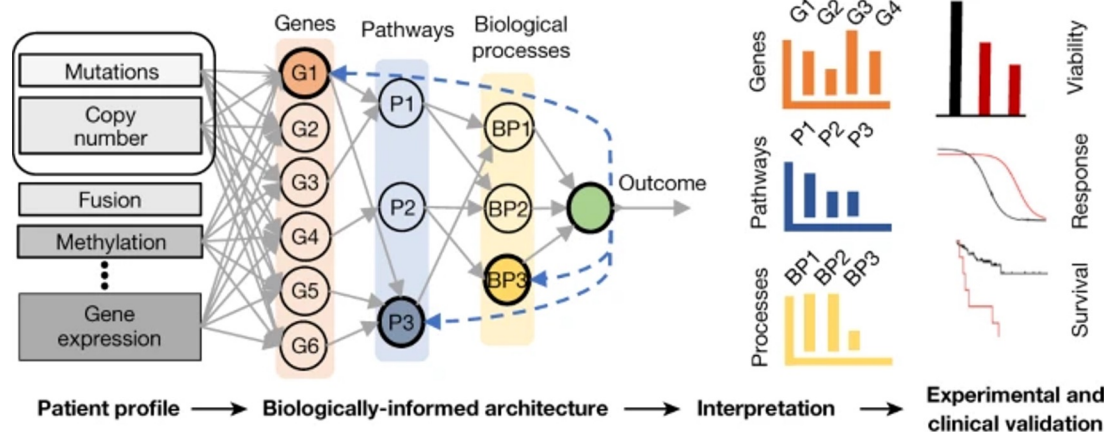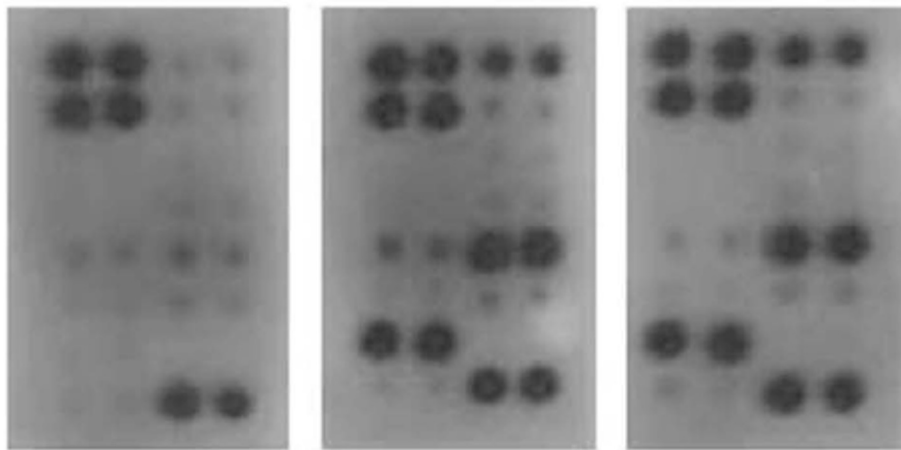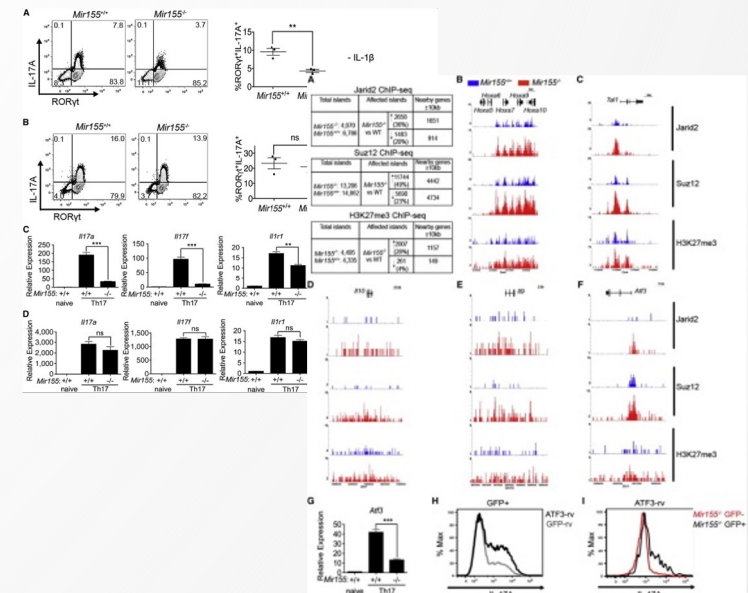# Science, Inference & Language

## Experiment /Observations



## Analysis



## Conclusions

*miR-155 Activates Cytokine Gene Expression in Th17 Cells by Regulating the DNA-Binding Protein Jarid2 to Relieve Polycomb-Mediated Repression.*

"miR-155 Activates Cytokine Gene Expression in Th17 Cells by Regulating the DNA-Binding Protein Jarid2 to Relieve Polycomb-Mediated Repression."

| | Patients with SARS-Cov-2 confirmed by PCR | Patients without SARS-Cov-2 confirmed by PCR |
|---|---|---|
| Median age (IQR)—years | 63 (53–72) | 60 (49–73) |
| Male | 787/1,309 (60.1%) | 90/167 (53.9%) |
| Race/ethnicity—Hispanic | 577/1,268 (45.5%) | 62/167 (37.1%) |
| Race/ethnicity—African American | 278/1,268 (21.9%) | 46/167 (27.5%) |
| Race/ethnicity—White | 277/1,268 (21.8%) | 43/167 (25.7%) |
| Race/ethnicity—Asian | 73/1,268 (5.8%) | 5/167 (3.0%) |
| Race/ethnicity—Other | 63/1,268 (5.0%) | 11/167 (6.6%) |
| Obesity (BMI ≥30) | 465/1,176 (39.5%) | 34/149 (22.8%)[a] |
| Comorbidities—hypertension | 420/1,268 (33.1%) | 67/167 (40.1%) |
| Comorbidities—diabetes | 293/1,268 (23.1%) | 34/167 (20.4%) |
| Comorbidities—CKD | 167/1,268 (13.2%) | 27/167 (16.2%) |
| … | … | … |

**Del Valle et al. , *Nature Medicine* (2020)**

$$\frac{dx_1(t)}{dt} = x_2(t)$$
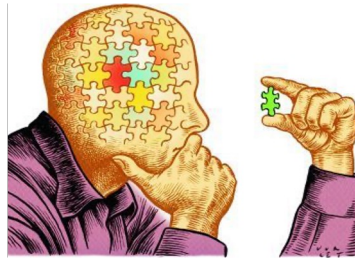
$$\frac{dx_2(t)}{dt} = ax_1(t) - bx_2(t)$$

$$\frac{d^2x_1(t)}{dt^2} = \frac{dx_2(t)}{dt}$$

where $x_1(t)$ is the serum concentration of cytokine and its rate of change by $x_2(t)$

Language & Abstraction!

# What if we could infer over scientific facts at scale?

## Reassemble & repurpose



Hypotheses Questions

Abductive Natural Language Inference (NLI)

Accumulated Knowledge

ANLI Models

# Abductive Reasoning

- First introduced by Peirce (1903).
- Inference to the best explanation.
- "Abduction is the mechanism via which we generate hypotheses about what we observe."
- Dialogues closely with assumed background knowledge.

Veen, Creative leaps in theory: the might of abduction (2021)

# **Abductive** Natural Language Inference (ANLI)

Inference to the best explanation
(facts, evidence)

**Claim:** Specialized cells protect the human body from disease-causing microbes by producing chemicals that destroy the microbes.

**True** | False

Multi-hop
Multi-premise

Why? (Explanation)

| Specialized cells are a source of chemicals that destroy disease-causing microbes. | disease-causing microbes have a negative impact on the body. |

~10.000 facts

# **Abductive** Natural Language Inference (ANLI)

Inference to the best explanation
(facts, evidence)

**Claim:** Specialized cells protect the human body from disease-causing microbes by producing chemicals that destroy the microbes.

**True** | False

Why? (Explanation)

Multi-hop
Multi-premise

Specialized cells are a source of chemicals that destroy disease-causing microbes.

disease-causing microbes have a negative impact on the body.

Encoding scientific statements

~10.000 facts

# **Abductive** Natural Language Inference (ANLI)

Inference to the best explanation
(facts, evidence)

**Claim:** Specialized cells protect the human body from disease-causing microbes by producing chemicals that destroy the microbes.
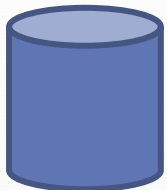
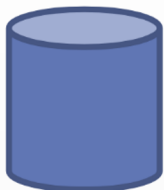**True** | False

Why? (Explanation)

Multi-hop
Multi-premise

Specialized cells are a source of chemicals that destroy disease-causing microbes.

disease-causing microbes have a negative impact on the body.

Encoding inference relations

~10.000 facts

# Expert-level scientific inference & explanation

**<u>Claim:</u>** BRCA2 promotes the joining of undamaged homologous repair molecules via RAD51 homolog 1 in humans.

BRCA2 and RAD51 homolog 1 are both involved in HRR in humans.

The binding of BRCA2 and RAD51 homolog 1 catalyzes the joining of undamaged homologous molecules.

RAD51 is a eukaryotic gene that encodes the RAD51 homolog gene.

BRCA2 promotes the assembly of RAD51 homolog 1 onto SS DNA in HRR.

BRCA2 is a human protein involved in DSB DNA break repair via HRR

BRCA2 is a human protein involved in HRR.

HRR is a DSB DNA repair process wherein damaged DNA is replaced by undamaged homologous molecules from sister chromatids or paternal/maternal copies of chromosomes.

BRCA2 is a human gene that encodes the BRCA2 protein.

BRCA2 protein is a tumour suppressor involved in HRR.

HRR repairs damage to DNA using information copied from a homologous undamaged molecule.

HRR is the primary process for repairing DNA double strand breaks.

**~1.000.000.000 facts**

Undamaged homologous molecules are provided by sister chromatids or paternal/maternal copies of chromosomes.

# Prostate cancer patient with loss of BRCA2 may benefit from PARP1 inhibition

Patients with loss of BRCA2 may benefit from PARP1 inhibition due to synthetic lethality causing cells to rely on a singular mechanism to repair cumulative damage to DNA.

PARP inhibitors cause replication-associated DSBs by preventing SS break repair, relying on defective HRR and error prone NHEJ to repair DNA.

Synthetic lethality is when co-occurrence of multiple genetic events results in cell death.

Loss of BRCA2 causes chromosome breakage.

BRCA2 is a human protein involved in HRR.

BRCA2 is a human gene that encodes the BRCA2 protein.

BRCA2 protein is a tumour suppressor that is involved in chromosomal stability.

Inhibition of PARP results in collapsed replication forks and DSB.

Inhibiting PARP results in accumulation of SS breaks.

SS breaks collapse replication forks and trigger HRR.

Inhibiting PARP results in accumulation of SS breaks.

PARP1 is involved in the recognition and repair of DNA damage in SS DNA damage repair.

PARP1 synthesis PAR which recruits repair proteins to sites of DNA damage

PARP1 detects and binds to sites of SS DNA damage.

PARP1 synthesises PAR.

PAR recruits repair proteins to damaged DNA site.

Loss of BRCA2 drives can development via genomic inst

Loss of BRCA2 may cause increased genomic instability.

Increas is a

Loss of BRCA2 causes the cell to default to NHEJ repair processes.

NHEJ does not use a ter DSB and can cause incr instability.

Loss of BRCA2 prevents the joining of undamaged repair molecules in HRR

In the absence of functional HRR genes, DNA repair defaults to NHEJ.

BRCA2 promotes the joining of undamaged homologous repair molecules via RAD51 homolog 1 in humans.

BRCA2 and RAD51 homolog 1 are both involved in HRR in humans.

The binding of BRCA2 and RAD51 homolog 1 catalyzes the joining of undamaged homologous molecules.

RAD51 is a eukaryotic gene that encodes the RAD51 homolog gene.

BRCA2 promotes the assembly of RAD51 homolog 1 onto SS DNA in HRR.

BRCA2 is a human protein involved in DSB DNA break repair via HRR

BRCA2 is a human protein involved in HRR.

HRR is a DSB DNA undamaged pat

BRCA2 is a human gene that encodes the BRCA2 protein.

BRCA2 protein is a tumour suppressor involved in HRR.

HRR is the primary process for repairing DNA double strand breaks.

## Key
In vivo
In vitro
Clinical Trials
Case series
Standard practice
External curated database
External uncurated database
Start/ End argumentation

Weak evidence ············
Good evidence – – – – – – –
Strong evidence ————

Patients living in the San Francisco area with ErbB2+ breast cancer, a body weight > 60 kg, and a history of treatment with Cyclophosphamide in the last year, are eligible for this clinical trial.
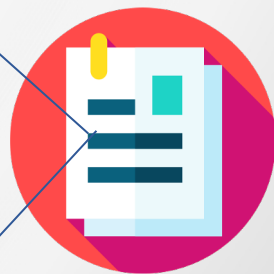


## Clinical Trial Report - Eligibility Criteria

**Inclusion criteria**
- Patients with a history of chemotherapy treatment within the last 24 months.
- Age ≥ 60 years
- HER2-positive T1 histologically confirmed invasive carcinoma of the breast.
- Body weight > 110 lbs
- Patients be California residents

**Exclusion criteria**
- Pregnant women

# The Neural Perspective: Language Models

- Probability distributions over strings of text.

The students opened their …
The students opened their <u>books</u>

(predicted)

**S** = The students opened their books

**P(S)** = P(The) x P(students | The) x P(opened | The students) x P(their | The students opened) x P(books | The students opened their)

# Neural Language Models



output distribution

$$\hat{y} = \text{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

concatenated word embeddings

$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hot vectors

$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

Kapronczay, Towards Data Science (2021)

# Transformers

1. Positional Encodings
2. (Multi-head) Self-Attention
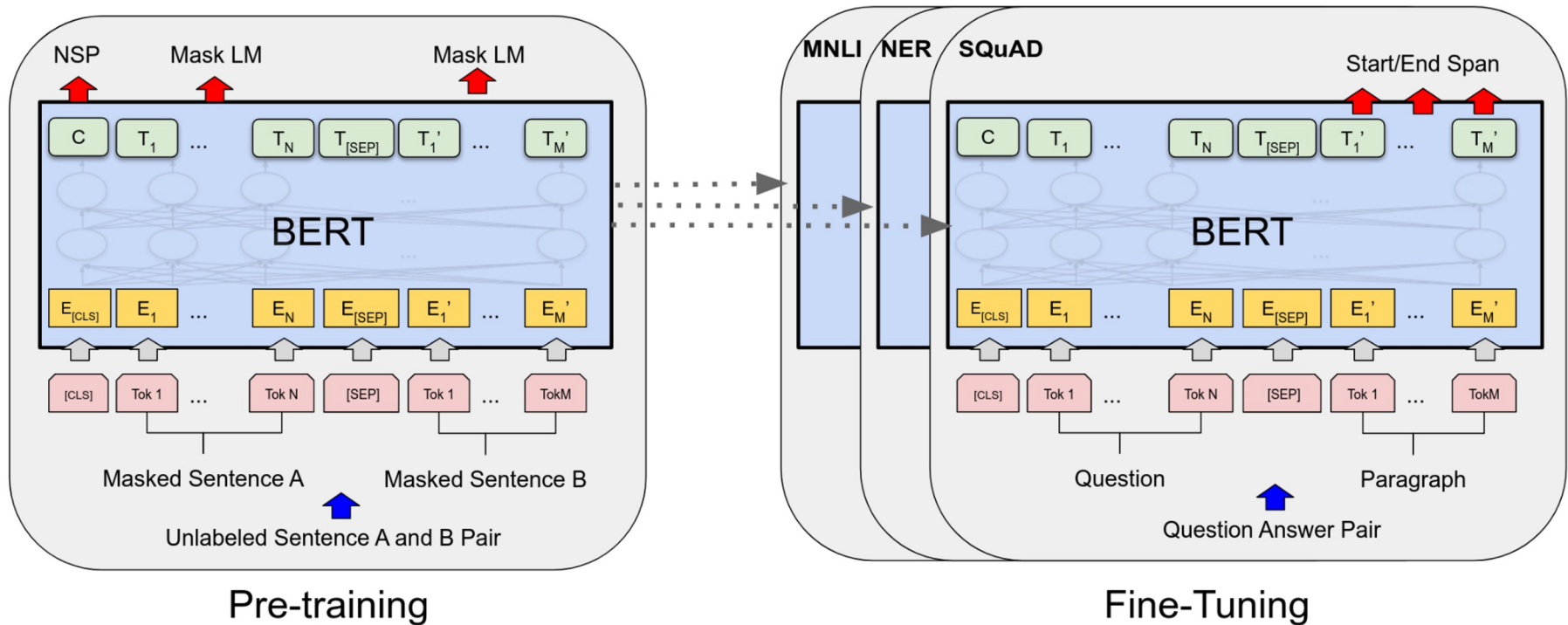
Vaswani et al, NeurIPS (2017)

# BERT: Bidirectional Encoder Representations

## from Transformers

Self-attention allows a a model to assign a meaning to a term in a complex context.



Devlin, Chang, Lee, Toutanova, CoRR (2018)

# Trust me, I am a Language Model.

Here is a sequence for a protein:

[START_AMINO]MEEPQSDPSVEPPLSQETFSDLWKLLPE...[END_AMINO]

And here is an isomeric SMILES for a compound:

[START_I_SMILES]CC(O)(P(=O)(O)O)P(=O)(O)O[END_I_SMILES]

**Question:** Will the the chemical compound be active against this protein?

**Answer:** No

# Trust me, I am a Language Model.

**Prompt**

The formula for Bessel's differential equation is:

**Generated Answer**

$$x^2 \frac{d^2 y}{dx^2} + x \frac{dy}{dx} + \left(x^2 - \alpha^2\right) y = 0$$

**Prompt**

Sulfuric acid reacts with sodium chloride, and gives _____ and _____:

\[ \ce{ NaCl + H2SO4 ->

**Generated Answer**

$$NaCl + H_2SO_4 \longrightarrow NaHSO_4 + HCl$$

# Why Meta's latest large language model survived only three days online

Galactica was supposed to help scientists. Instead, it mindlessly spat out biased and incorrect nonsense.

By Will Douglas Heaven                                    November 18, 2022

**Michael Black**
@Michael_J_Black · Follow

I asked #Galactica about some things I know about and I'm troubled. In all cases, it was wrong or biased but sounded right and authoritative. I think it's dangerous. Here are a few of my experiments and my analysis of my concerns. (1/9)
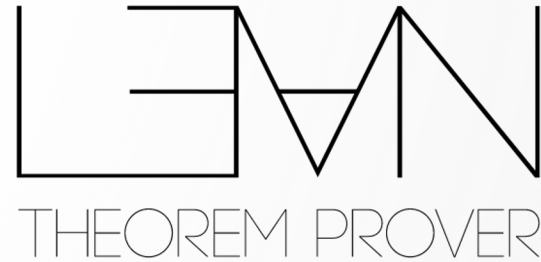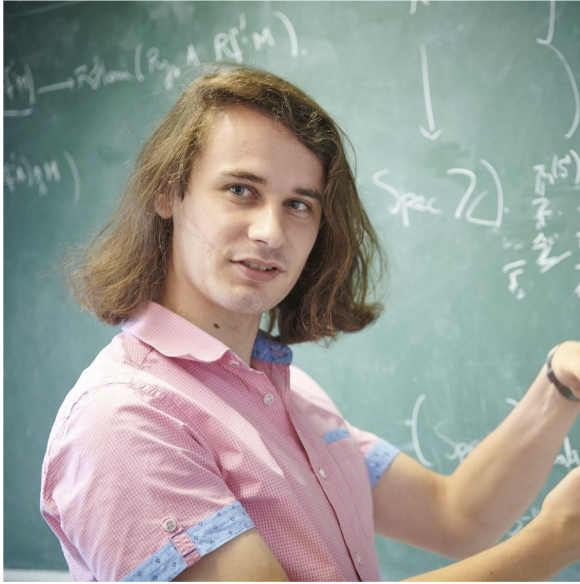
7:47 AM · Nov 17, 2022

Ouch!

**Julian Togelius**
@togelius · Follow

My considered opinion of Galactica: it's fun, impressive, and interesting in many ways. Great achievement. It's just unfortunate that it's being touted as a practical research tool, and even more unfortunate that it suggests you use it to write complete articles.
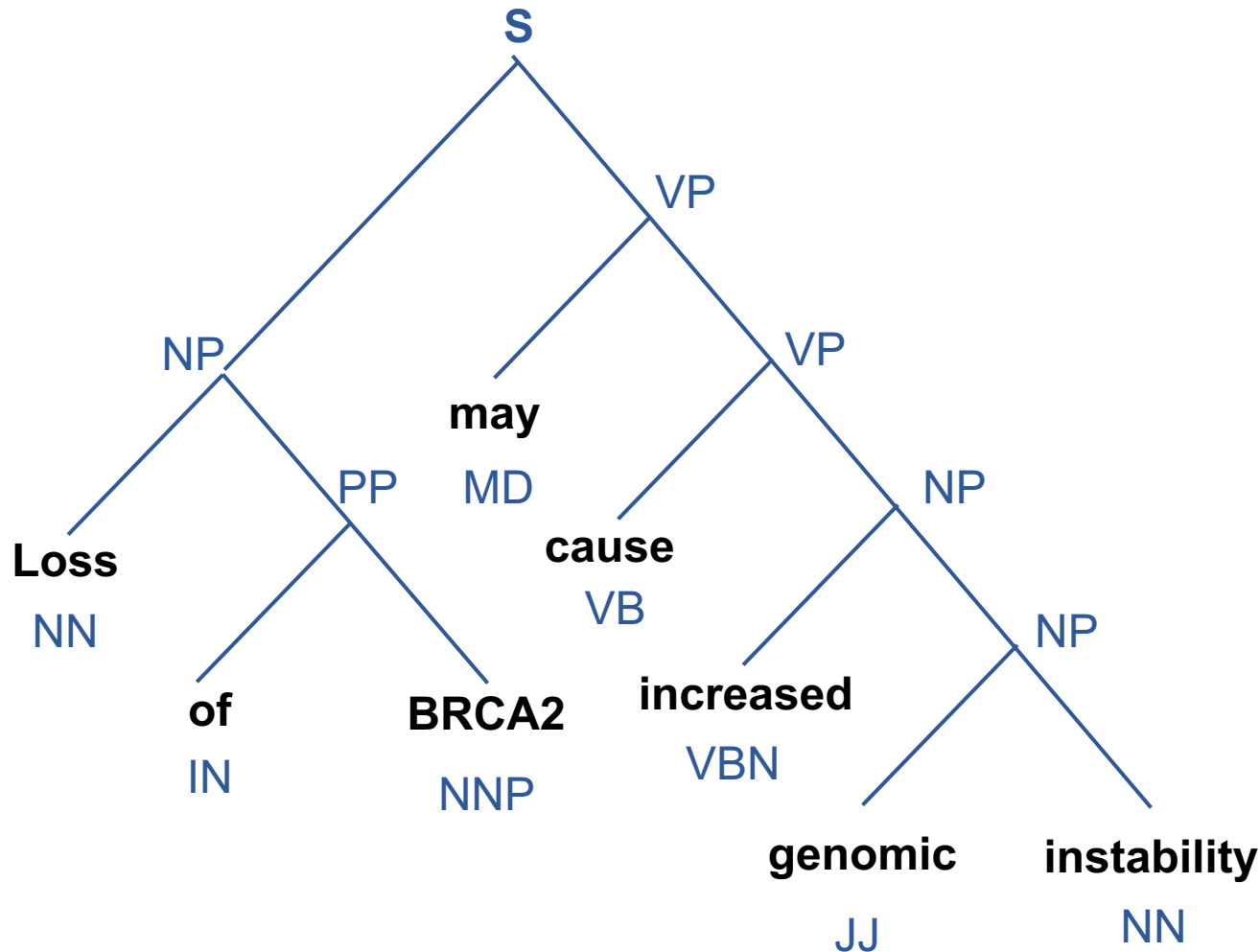
# The Liquid Tensor Experiment

## Why do I want a formalization?

— "with this theorem, the hope that the condensed formalism can be fruitfully applied to real functional analysis stands or falls. I think the theorem is of utmost foundational importance, **so being 99.9% sure is not enough**."
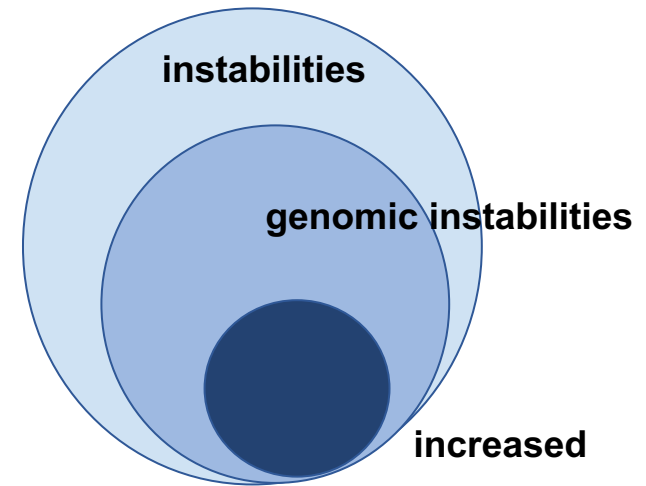
**Verifiability**

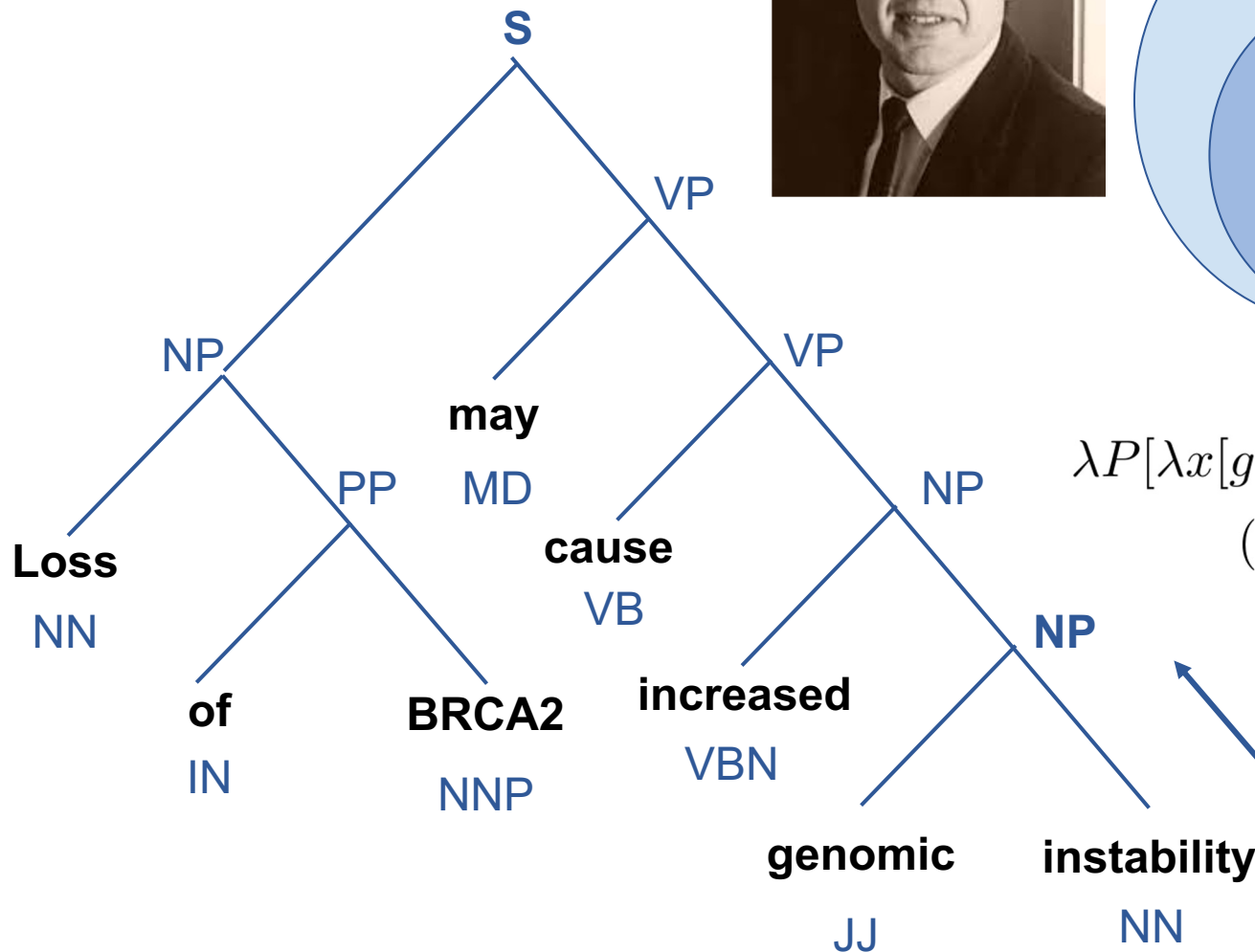https://xenaproject.wordpress.com/2020/12/05/liquid-tensor-experiment/

# The Formal Perspective

Loss of BRCA2 may cause increased genomic instability.

# The Formal Perspective



$\lambda P[\lambda x[genomic(x) \land P(x)]]$
$(\lambda y[instability(y)])$

# Scientific inference

**Scientific discourse**
- Step-wise explicit (verbalised) inference.
- Formal, verifiable argument & explanation.
- Preserving the positive aspects of LLMs.
- Improving control.

Large Language
Models (LLMs)

Formal

Neuro

**Neuro-symbolic**

Symbolic

# Scientific inference
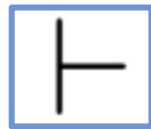
**Scientific discourse**
- Step-wise explicit (verbalised) inference.
- Formal, verifiable argument & explanation.
- Preserving the positive aspects of LLMs.
- Improving control.

$$\Gamma \vDash \Phi$$

$\Gamma$ semantically entails $\Phi$

$$\Gamma \vdash \Phi$$

$\Gamma$ proves $\Phi$

$\vdash$

- interpretability
- control (inference guarantees)

# Encoding scientific statements

# Semantic Role Labeling

Large Language Models ← Formal

Neuro                                              Symbolic

| Animals | require | food | for survival | . |
|---------|---------|------|--------------|---|
| ARG0    | V       | ARG1 | ARGM-PRP     |   |

- Lightweight representation (a little semantics goes a long way).
- Robust parsers.
- Expressive semantic roles.
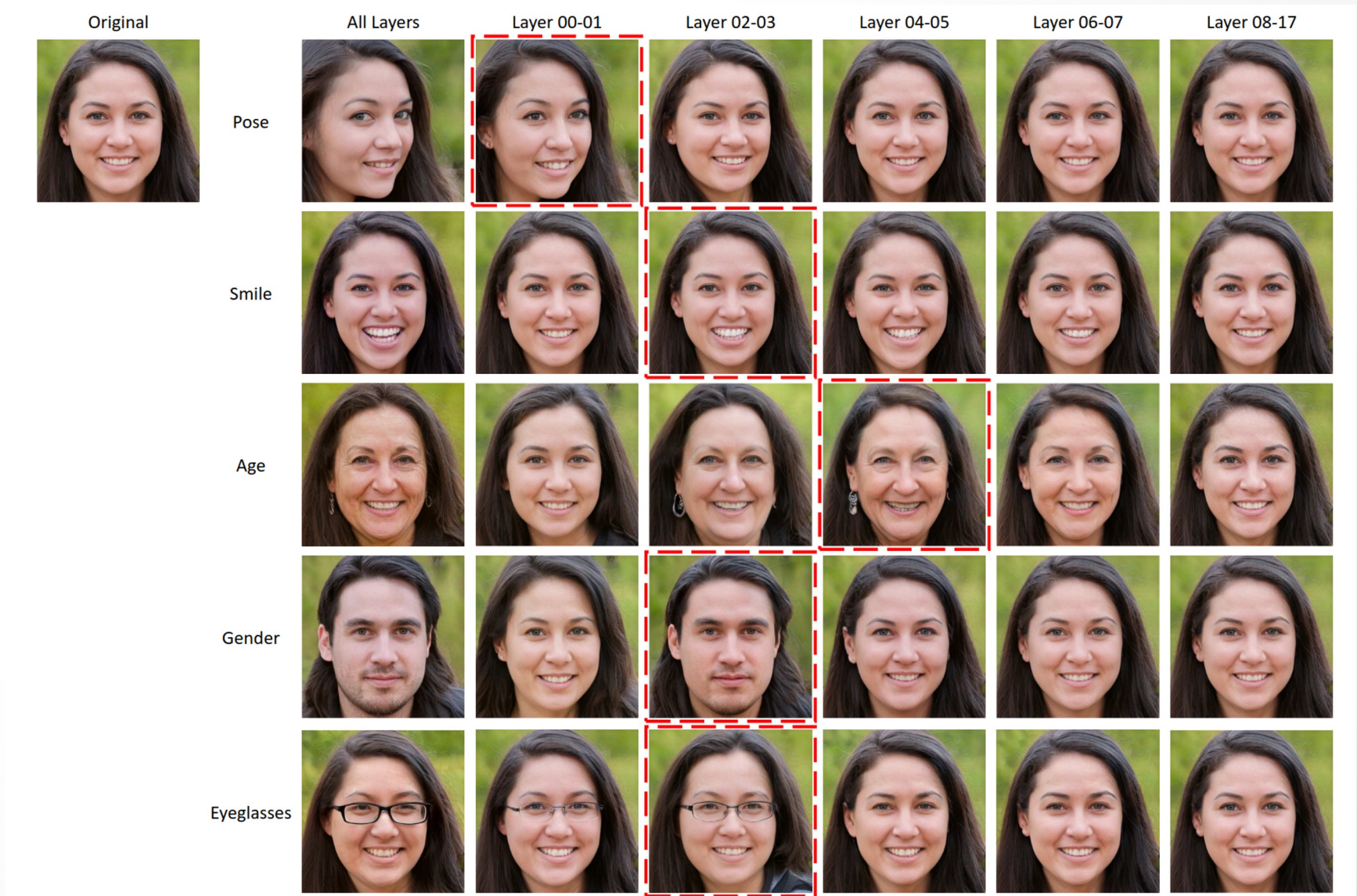
| | |
|---|---|
| ARGM-DIR | Directionals. E.g. all waves transmit energy **from one place to another** |
| ARGM-PNC | Purpose. E.g. many animals blend in with their environment **to not be seen by predators** |
| ARGM-CAU | Cause. E.g. cold environments sometimes are white in color **from being covered in snow** |
| ARGM-PRP | Purpose. E.g. a pot is made of metal **for cooking** |
| ARGM-EXT | Extent. E.g. as the amount of oxygen exposed to a fire increases the fire will burn **longer** |
| ARGM-LOC | Location. E.g. a solute can be dissolved **in a solvent** when they are combined |
| ARGM-MNR | Manner. E.g. fast means **quickly** |
| ARGM-MOD | Modal verbs. E.g. atom **can** not be divided into smaller substances |
| ARGM-DIS | Discourse. E.g. if something required by an organism is depleted **then** that organism must replenish that something |
| ARGM-GOL | Goal. E.g. We flew **to Chicago** |
| ARGM-NEG | Negation. E.g. cactus wrens building nests in cholla cacti does **not** harm the cholla cacti |
| ARGM-ADV | Adverbials |
| ARGM-PRD | Markers of secondary predication. E.g. |
| ARGM-TMP | Temporals. E.g. a predator **usually** kills its prey to eat it |
| O | Empty tag. |
| V | Verb. |
| ARG0 | Agent or Causer. E.g. **rabbits** eat plants |
| ARG1 | Patient or Theme. E.g. rabbits eat **plants** |
| ARG2 | indirect object / beneficiary / instrument / attribute / end state. E.g. animals are **organisms** |
| ARG3 | start point / beneficiary / instrument / attribute. E.g. sleeping bags are designed **to keep people warm** |
| ARG4 | end point. E.g. when water falls from the sky that |

# Generative Models

**Disentanglement**

**Z**

**VAE**

(Loss of BRCA2) **causes** (the cell) to default to (NHEJ repair processes).

(Loss of BRCA2) **causes** (the cell) to default to (NHEJ repair processes).

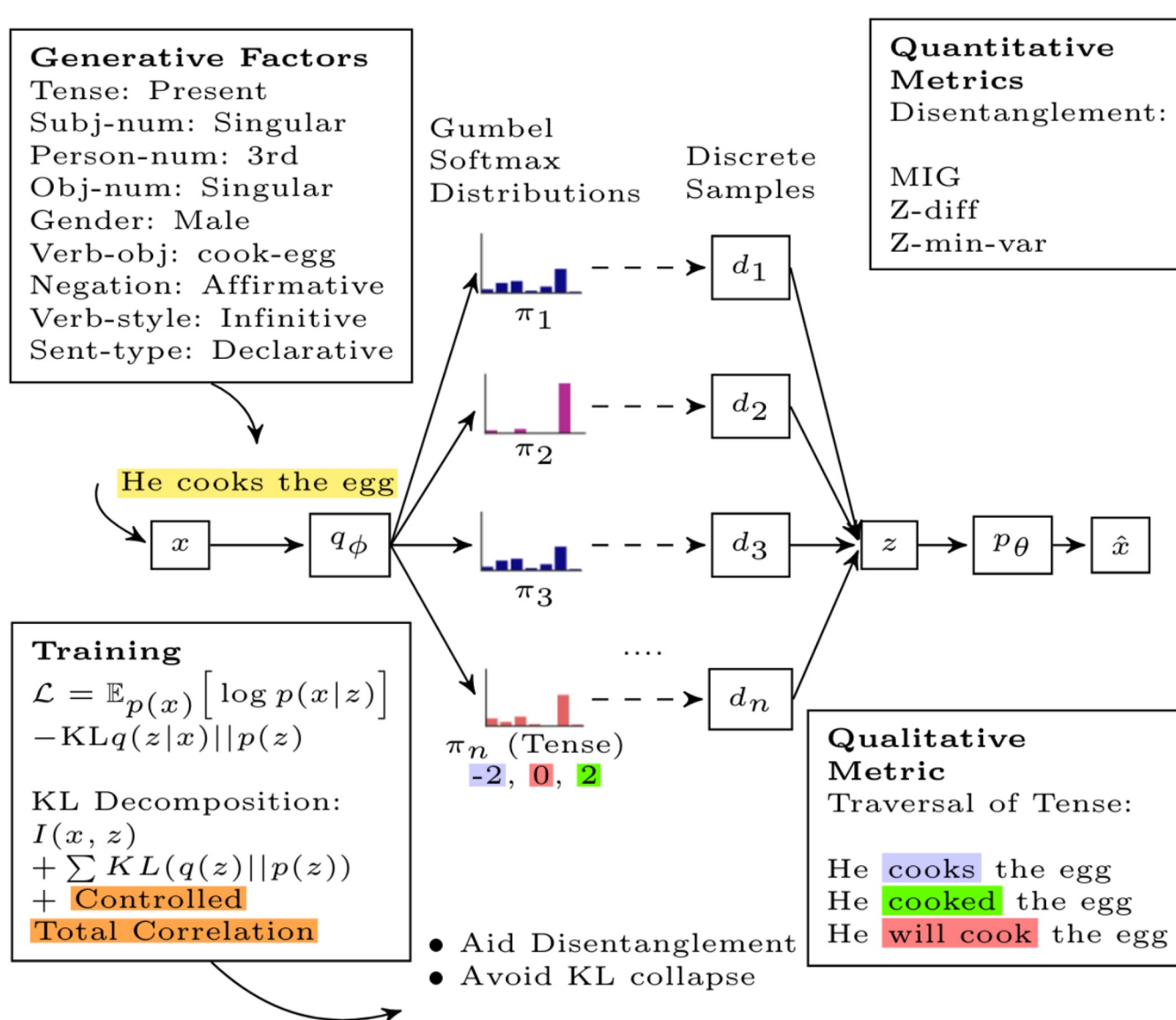**Integrating Syntactic and Semantic Structure into the latent space**

# Generative Models



*InterFaceGAN: Interpreting the Disentangled Face Representation Learned by GANs*

# Syntactic Disentanglement



**Generative Factors**
Tense: Present
Subj-num: Singular
Person-num: 3rd
Obj-num: Singular
Gender: Male
Verb-obj: cook-egg
Negation: Affirmative
Verb-style: Infinitive
Sent-type: Declarative

He cooks the egg

Gumbel
Softmax
Distributions

Discrete
Samples

$\pi_1$

$\pi_2$

$\pi_3$

$\pi_n$ (Tense)
-2, 0, 2

$d_1$

$d_2$

$d_3$

$d_n$

$x$ $\longrightarrow$ $q_\phi$ $\longrightarrow$ $z$ $\longrightarrow$ $p_\theta$ $\longrightarrow$ $\hat{x}$

**Quantitative
Metrics**
Disentanglement:

MIG
Z-diff
Z-min-var

**Training**
$$\mathcal{L} = \mathbb{E}_{p(x)}\Big[\log p(x|z)\Big]$$
$$-\mathrm{KL}\, q(z|x)||p(z)$$

KL Decomposition:
$I(x, z)$
$+ \sum KL(q(z)||p(z))$
$+$ Controlled
Total Correlation

• Aid Disentanglement
• Avoid KL collapse

**Qualitative
Metric**
Traversal of Tense:

He cooks the egg
He cooked the egg
He will cook the egg

Mercatali & Freitas, EMNLP Findings (2021)

# Syntactic Disentanglement

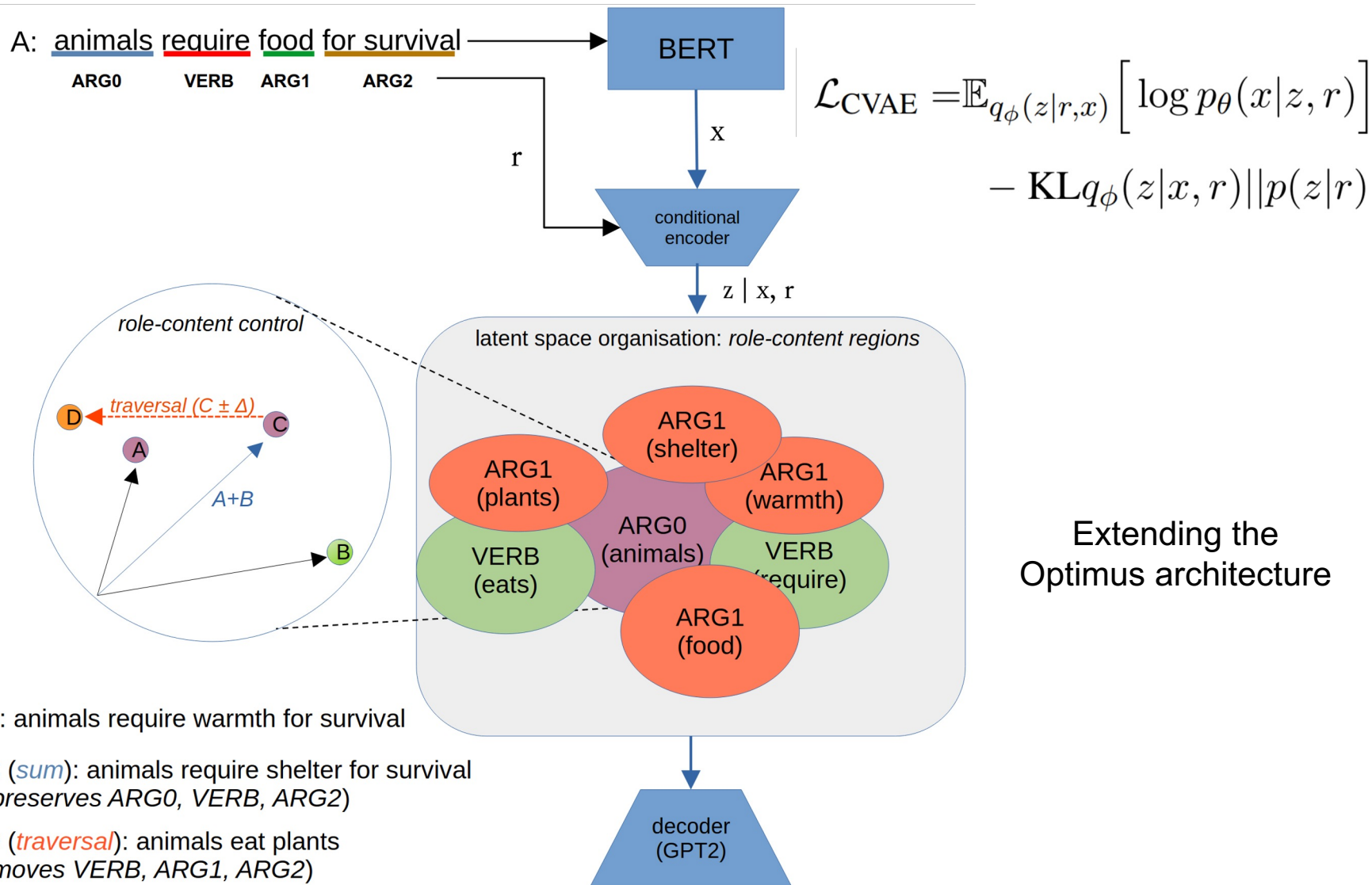(Loss of BRCA2) **causes** (the cell) to default to (NHEJ repair processes).
(Loss of BRCA2) **caused** (the cell) to default to (NHEJ repair processes).
(Loss of BRCA2) **does not cause** (the cell) to default to (NHEJ repair processes).

| | Tense | Subject-number |
|---|---|---|
| input | you will not attend the party | we will not attend the party |
| βVAE | you will not attend the party<br>you will not sign the paper<br>you will not attend the party | we will not attend the party<br>he will not attend the party |
| JointVAE | you will not attend the party<br>you did not join the wedding<br>you do not attend the party | we will not attend the party<br>you will not attend the party |
| DCTC | you will not attend the party<br>you did not attend the party<br>you do not attend the party | we will not attend the party<br>i will not attend the party |

**Latent traversal**

Mercatali & Freitas, EMNLP Findings (2021)

# Syntactic-Semantic Disentanglement



$$\mathcal{L}_{\text{CVAE}} = \mathbb{E}_{q_\phi(z|r,x)}\left[\log p_\theta(x|z,r)\right]$$
$$- \text{KL}\, q_\phi(z|x,r)||p(z|r)$$

Extending the
Optimus architecture

B: animals require warmth for survival

C (*sum*): animals require shelter for survival
(*preserves ARG0, VERB, ARG2*)

D (*traversal*): animals eat plants
(*moves VERB, ARG1, ARG2*)

Zhang, Carvalho, Pratt-Hartmann, Freitas, arXiv:2210.06230 (2022)

# Syntactic-Semantic Disentanglement

an automobile is a kind of vehicle

an automobile requires a driver to move it
an automobile is a kind of object

an airplane is a kind of vehicle
a car is a kind of vehicle

an airplane is used for carrying passengers
an airplane is a kind of object

**Latent traversal**

animals require food for survival
animals require warmth for survival

animals eat plants
animals produce milk
animals usually eat plants
animals eat berries ; plants
animals require food to survive
animals require shelter to survive
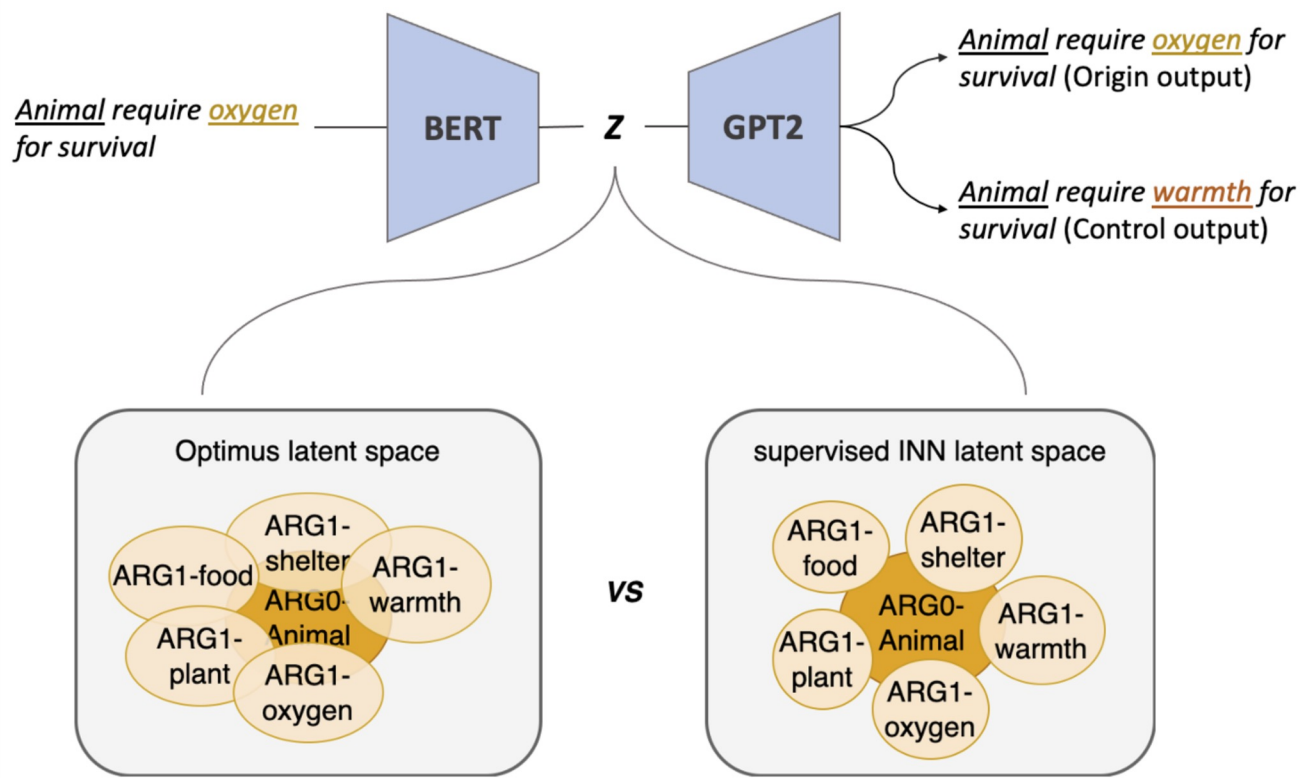animals adapt to changing environments
animals obtain spices for cooking

**Vector arithmetic
(addition)**

Zhang, Carvalho, Pratt-Hartmann, Freitas, arXiv:2210.06230 (2022)

# Improving Separability

**Adding a flow-based INN component to improve separability**

$$\mathcal{L}_{\text{sup}} = - \mathbb{E}_{x \sim p_{cluster}(x)} \frac{\left[T(E(x)) - \mu_{cluster}\right]^2}{1 - \sigma^2} - \log \left|T'(E(x))\right|$$

## Interpolation

**humans eat seeds**
1. humans eat fruits
2. humans eat seeds
3. humans eat insects
4. humans eat meat
5. humans eat plants
6. some animals eat prey
7. some animals must eat to survive
8. some animals must hunt for food
9. some animals must hunt their prey to survive
**some animals must hunt to survive**

**Interpolation**

**Latent traversal**

## Input: some animals must hunt to survive

dim01: **some animals** must hunt for food
dim01: **some animals** must hunt prey to survive
dim01: **some animals** need to hunt to survive

dim12: **an animal** needs to breathe to survive
dim12: **an animal** can fly without air
dim12: **a predator** must hunt to survive

# Data Augmentation

| Role-content | Augmented sentences |
| --- | --- |
| ARG0-animal | an animal requires energy to move <br> animals produce offspring <br> some adult animals lay eggs <br> an animal requires shelter <br> an animal can use its body to breathe |
| ARG0-human | humans travel sometimes <br> humans usually use gasoline <br> humans sometimes endanger themselves <br> humans use coal to make food <br> humans depend on pollinators for survival |
| PRED-are | wheels are a part of a car <br> lenses are a part of eyeglasses <br> toxic chemicals are poisonous <br> green plants are a source of food for animals <br> copper and zinc are two metals |
| PRED-mean | summit mean the top of the mountain <br> colder mean a decrease in heat energy <br> helping mean something can be done better <br> cleaner mean ( less ; lower ) in pollutants <br> friction mean the product of a physical change |

# Representing concepts and definitions

- Essential attributes of a conceptualisation.
- Abundance of NL definitions in scientific discourse.
- Definition RL: Decomposing conceptual components.

**DEFINIENDUM**   **DIFFERENTIA QUALITY** **SUPERTYPE** **DIFFERENTIA-EVENT**

**Homologous recombination repair** is a **DNA repair process** that **includes the invasion of an undamaged DNA molecule by a damaged molecule of identical or very similar sequence**.

**DSR Optimus** — **Interpolation**

a migratory aquatic bird found in the temperate regions
of the northern hemisphere
1 a migratory bird of the eastern Mediterranean
2 a marine gastropod of the subfamily
3 a terrestrial aquatic mammal of the family
4 a terrestrial aquatic mammal of the suborder
5 a terrestrial invertebrate
6 a microscopic organism or invertebrate
a microscopic terrestrial animal or protozoan

an automobile
1 a motorcycle
a bicycle

**Latent traversal**

**ADD**
a flying machine
a flying creature
a flying dinosaur
a flying robot
a flying object

**AVG**
to make four copies of
to make five copies of
to make one copy of
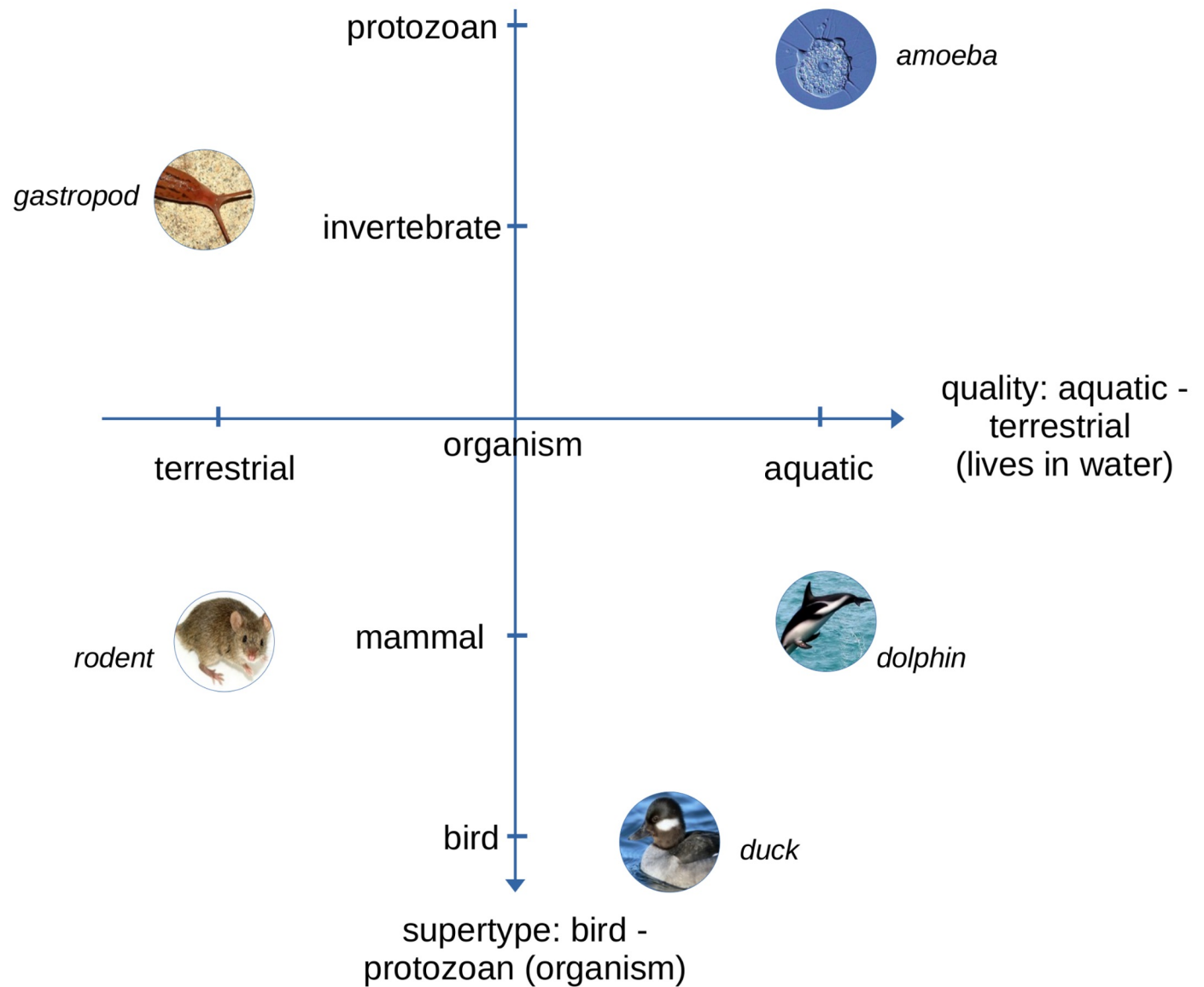to make two copies of
to make 3 copies of

**SUB**
a female monarch
a monarch
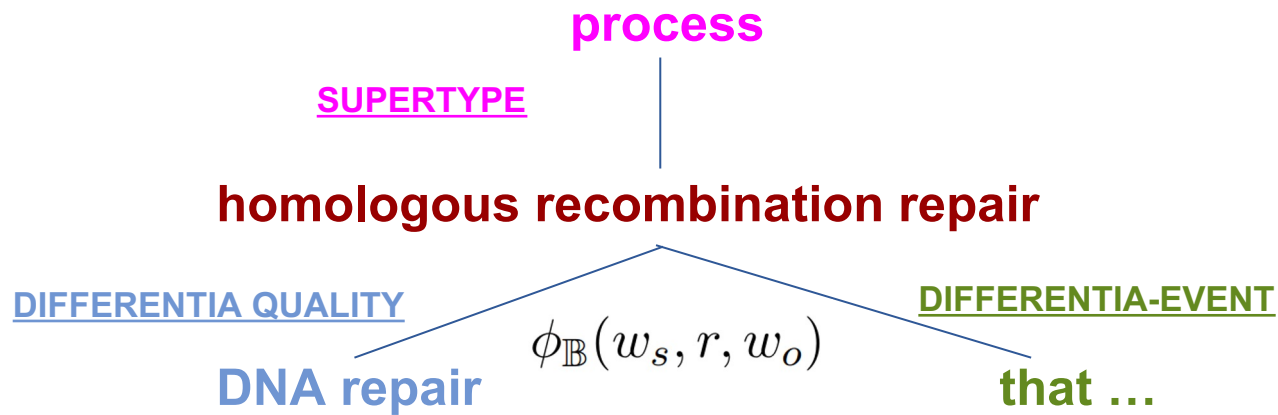the subnormal condition in females originating from...
the normal female pregnancy associated with some
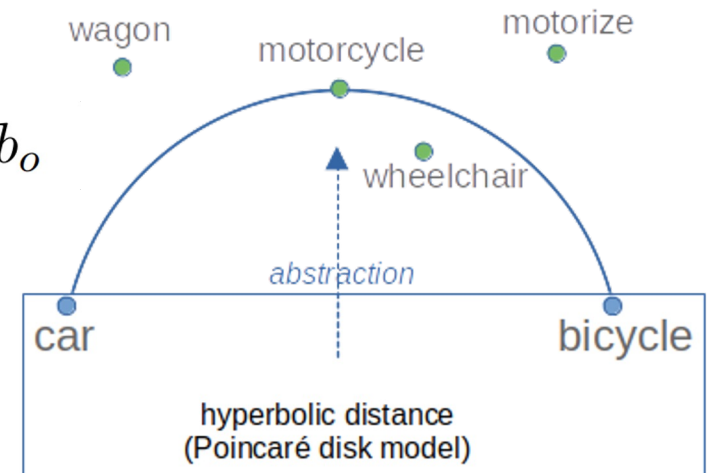the female given name in the Japanese game...

Carvalho, Mercatali, Zhang, Freitas, arXiv:2210.02898 (2022)

protozoan — amoeba

gastropod

invertebrate

quality: aquatic -
terrestrial
(lives in water)

organism

terrestrial

aquatic

mammal — dolphin

rodent

bird — duck

supertype: bird -
protozoan (organism)

Carvalho, Mercatali, Zhang, Freitas, arXiv:2210.02898 (2022)

# Multi-relational Hyperbolic Embeddings

Induction of a hierarchical, multi-relational, multi-resolution conceptual representation. Abstracts OOV words.

**process**

**SUPERTYPE**

**homologous recombination repair**

**DIFFERENTIA QUALITY**          **DIFFERENTIA-EVENT**

$\phi_{\mathbb{B}}(w_s, r, w_o)$

**DNA repair**          **that …**

Learn as a link prediction problem via a translational objective function in hyperbolic space.

$$\phi_{\mathbb{B}}(w_s, r, w_o) = -d_{\mathbb{B}}(\mathbf{h}_s^{(r)}, \mathbf{h}_o^{(r)})^2 + b_s + b_o$$

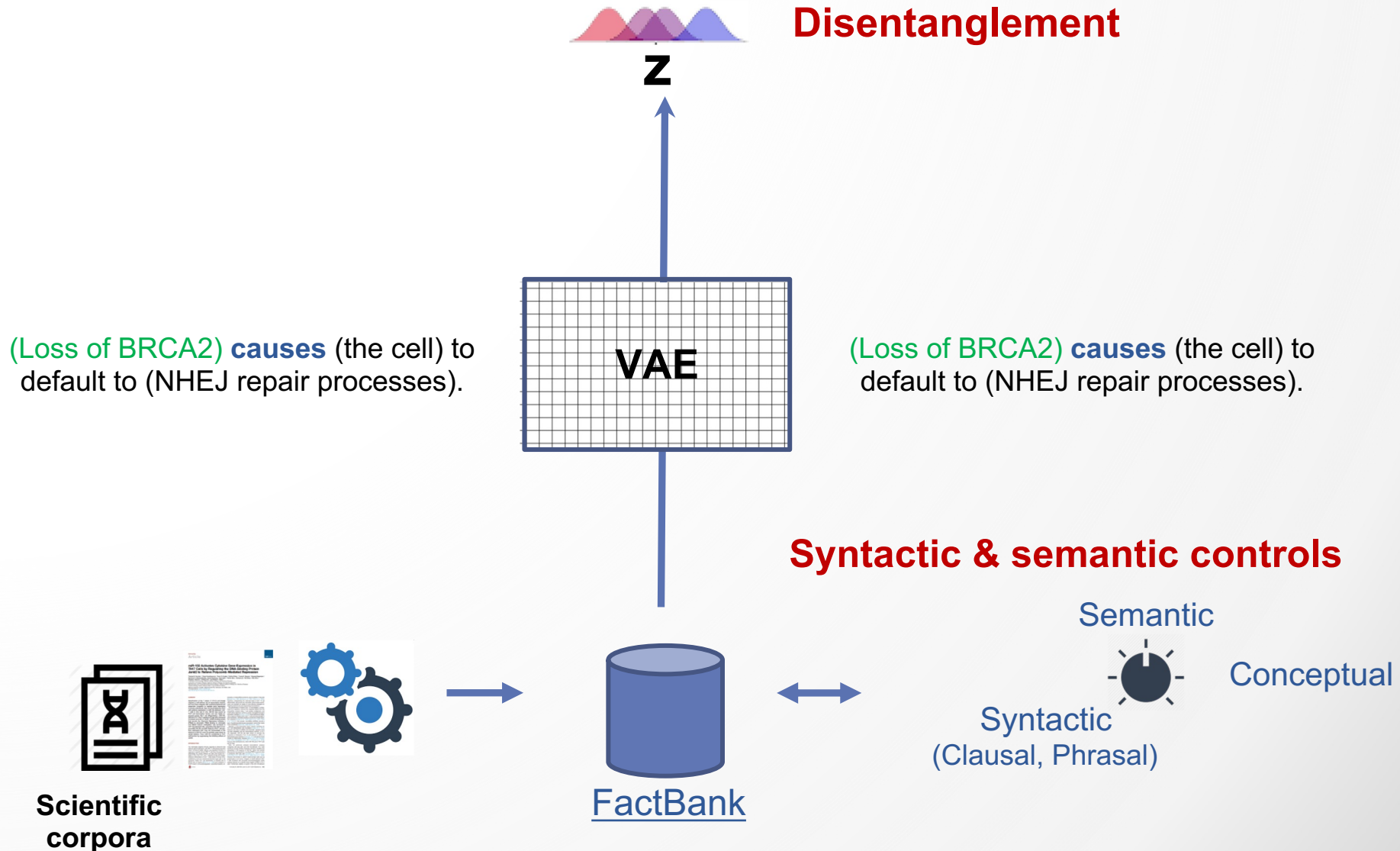$$d_{\mathbb{B}}(x, y) = \frac{2}{\sqrt{c}} tanh^{-1}(\sqrt{c}\|-x \oplus y\|)$$

wagon    motorcycle    motorize

wheelchair

abstraction

car    bicycle

hyperbolic distance
(Poincaré disk model)

!

| Model | SV-d | MEN-d | SV-t | MEN-t | SL999 | SCWS | 353 | RG | MT |
|---|---|---|---|---|---|---|---|---|---|
| **Transformers** | | | | | | | | | |
| SBERT (bert-base) | 13.5 | 27.8 | 13.3 | 30.6 | 15.1 | 37.8 | 20.0 | 68.1 | 22.3 |
| SBERT (bert-large) | 16.1 | 23.4 | 14.4 | 26.8 | 13.4 | 35.7 | 19.8 | 60.7 | 19.1 |
| SBERT (distilroberta) | 35.8 | 61.2 | 36.7 | 62.2 | 43.4 | 57.1 | 52.0 | 77.4 | 46.2 |
| SBERT (mpnet-base) | 45.9 | **64.9** | 42.5 | **67.5** | 49.5 | 58.6 | 56.5 | 81.3 | 45.3 |
| SBERT (t5-large) | **49.4** | 63.1 | **50.2** | 66.3 | **57.3** | 56.1 | 51.8 | **85.3** | 38.1 |
| **Multi-Relational** | | | | | | | | | |
| Euclidean ($d = 40$) | 39.1 | 62.9 | 35.7 | 65.4 | 36.3 | 58.2 | 52.1 | 80.9 | 45.0 |
| Euclidean ($d = 80$) | 44.1 | 65.6 | 39.5 | 66.2 | 41.2 | 58.4 | 55.8 | 78.0 | 42.4 |
| Euclidean ($d = 200$) | 47.3 | 67.0 | 41.0 | 67.6 | 43.4 | 60.6 | 55.4 | 78.1 | 44.6 |
| Euclidean ($d = 300$) | 47.9 | 68.3 | 43.1 | 69.1 | **44.7** | 61.0 | 54.4 | 79.0 | 46.0 |
| Hyperbolic ($d = 40$) | 36.7 | 66.2 | 34.3 | 66.4 | 31.8 | 58.5 | 50.6 | 75.5 | 52.7 |
| Hyperbolic ($d = 80$) | 42.7 | 68.2 | 40.7 | 68.6 | 38.3 | 61.4 | 59.2 | 81.0 | **59.1** |
| Hyperbolic ($d = 200$) | 48.8 | 71.8 | 44.7 | 73.2 | 40.7 | 63.5 | 64.9 | **81.6** | 57.6 |
| Hyperbolic ($d = 300$) | **50.6** | **72.6** | **45.4** | **74.2** | 42.3 | **63.9** | **66.3** | 80.5 | 56.1 |

T5-large: "Colossal Clean Crawled Corpus" (C4): ~750GB.
MR–Hyperbolic: "CPAE Dictionary": ~19MB.

# Sentence-level representation



**Disentanglement**

**z**

(Loss of BRCA2) **causes** (the cell) to default to (NHEJ repair processes).

**VAE**

(Loss of BRCA2) **causes** (the cell) to default to (NHEJ repair processes).

**Syntactic & semantic controls**

Semantic

Conceptual

Syntactic
(Clausal, Phrasal)

**Scientific corpora**

FactBank

# Encoding step-wise (multi-hop) inference

# Encoding Abductive (Explanatory) Reasoning

**h: <u>Shale</u> is a <u>sedimentary rock</u> that can be metamorphosed into <u>slate</u> by <u>increased pressure</u>.**

'<u>shale</u> is a kind of <u>sedimentary rock</u>'          '<u>high</u> is similar to <u>increase</u>'

'<u>extreme</u> means very <u>high</u> in value'

'<u>slate</u> is a type of <u>metamorphic rock</u>'

'exposure to <u>extreme</u> heat and <u>pressure</u> changes <u>sedimentary</u> and igneous <u>rock</u> into <u>metamorphic rock</u>'

**Abstraction, grounding**

**Abstraction**

Concrete facts tend to share key concepts with the hypotheses and can therefore be effectively retrieved by lexical relevance.

**h: Shale is a sedimentary rock that can be metamorphosed into slate by increased pressure.**

'shale is a kind of sedimentary rock'                     'high is similar to increase'

'extreme means very high in value'

'slate is a type of metamorphic rock'

'exposure to extreme heat and pressure changes sedimentary and igneous rock into metamorphic rock'

**Unification**

**Abstraction**

More uiversal scientific statements tend to be abstract and therefore difficult to rank by means of shared concepts.

**h: <u>Shale</u> is a <u>sedimentary rock</u> that can be metamorphosed into <u>slate</u> by <u>increased pressure</u>.**

'<u>shale</u> is a kind of <u>sedimentary rock</u>'         '<u>high</u> is similar to <u>increase</u>'

'<u>extreme</u> means very <u>high</u> in value'

'<u>slate</u> is a type of <u>metamorphic rock</u>'

'exposure to <u>extreme</u> heat and <u>pressure</u> changes <u>sedimentary</u> and igneous <u>rock</u> into <u>metamorphic rock</u>'

**Abstraction**

Proposes the composition of two scoring functions:

• A **<u>Relevance Score (RS)</u>** that represents the lexical relevance of a given fact.

• A **<u>Unification Score (US)</u>** that models the explanatory power of a fact according to its frequency in explanations for similar questions

Valentino, Thayaparan, Freitas, EACL (2021)

**Question(Q):**

What is an example of force producing heat?

**Candidate Answer (C₁):**

Two sticks getting warm when rubbed together

**Hypothesis (H₁):**

Two sticks getting warm when rubbed together is an example of force producing heat

**Grounding Facts:**

[✓] a stick is an object: $F_{G1}$

[✓] friction is a force: $F_{G2}$

[✗] a pull is a force: $F_{G3}$

[✓] to rub together means to move against: $F_{G4}$
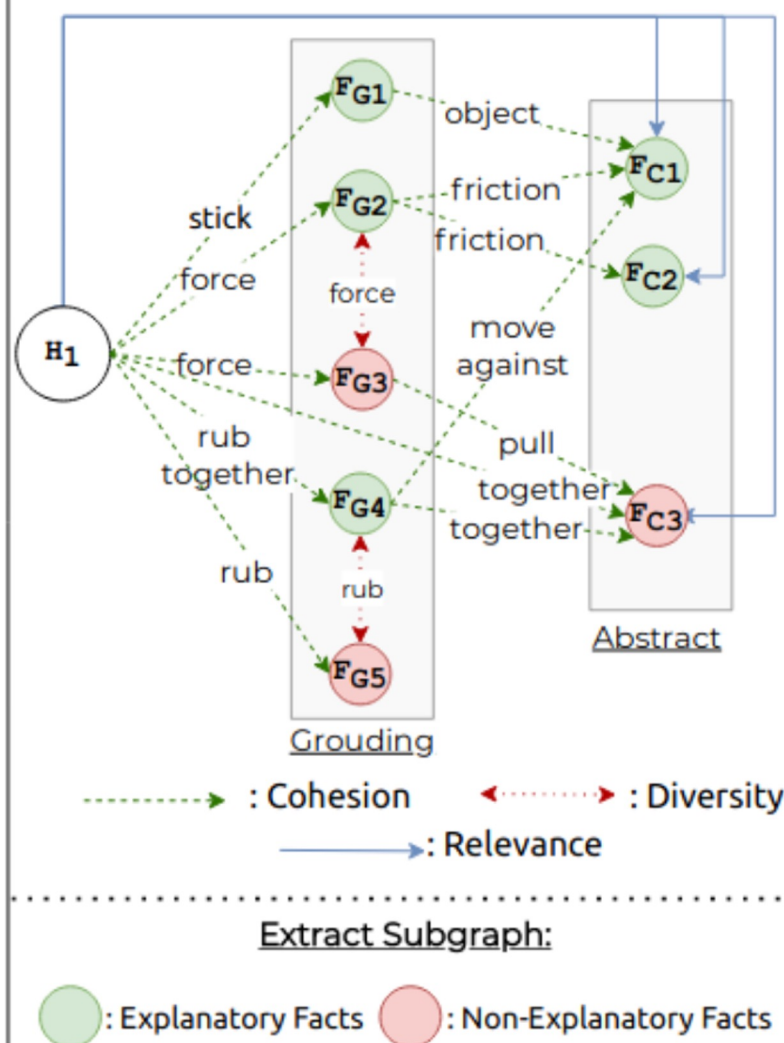
[✗] rubbing against something is kind of movement: $F_{G5}$

**Abstract Facts:**

[✓] friction occurs when two object's surfaces move against each other: $F_{C1}$

[✓] friction causes the temperature of an object to increases: $F_{C2}$

[✗] magnetic attraction pulls two objects together: $F_{C3}$

[✓]: Explanatory Facts
[✗]: Non-Explanatory Facts

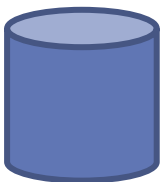Thayaparan, Valentino, Freitas, ACL Findings (2021)

**Question(Q):**

What is an example of force producing heat?

**Candidate Answer (C₁):**

Two sticks getting warm when rubbed together

**Hypothesis (H₁):**

Two sticks getting warm when rubbed together is an example of force producing heat

**Grounding Facts:**

[✓] a stick is an object: $F_{G1}$
[✓] friction is a force: $F_{G2}$
[✗] a pull is a force: $F_{G3}$
[✓] to rub together means to move against: $F_{G4}$
[✗] rubbing against something is kind of movement: $F_{G5}$

**Abstract Facts:**

[✓] friction occurs when two object's surfaces move against each other: $F_{C1}$
[✓] friction causes the temperature of an object to increases: $F_{C2}$
[✗] magnetic attraction pulls two objects together: $F_{C3}$

[✓]: Explanatory Facts
[✗]: Non-Explanatory Facts

**For each Candidate Hypothesis:**

Fact Graph Construction:

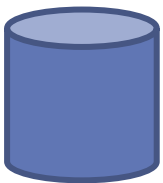object
friction
friction
move against
pull
together
together
stick
force
force
force
rub together
rub
rub

$H_1$  $F_{G1}$ $F_{G2}$ $F_{G3}$ $F_{G4}$ $F_{G5}$  $F_{C1}$ $F_{C2}$ $F_{C3}$

Grounding

Abstract

- - - - → : Cohesion      ◄·······► : Diversity

————→ : Relevance

Extract Subgraph:

◯ : Explanatory Facts   ◯ : Non-Explanatory Facts

Thayaparan, Valentino, Freitas, ACL Findings (2021)

Abstract Facts

Grounding Facts

$h_i$

$f_k$

...

$\{f_1, \ f_2, \ f_3, \ \cdots, \ f_k\}$

$W_{ik}$

$G \ = \ (H, \ F, \ E, \ W)$

Retrieve relevant facts

Construct a weighted fact graph

Extract subgraph via ILP optimization

**Relevance**

**Diversity**

$$D(f_j^{h_i}, f_k^{h_i}) = -1 \frac{|t_{h_i}(f_j^{h_i}) \cap t_{h_i}(f_k^{h_i})|}{max(|t_{h_i}(f_j^{h_i})|, |t_{h_i}(f_k^{h_i})|)}$$

**Saturation**

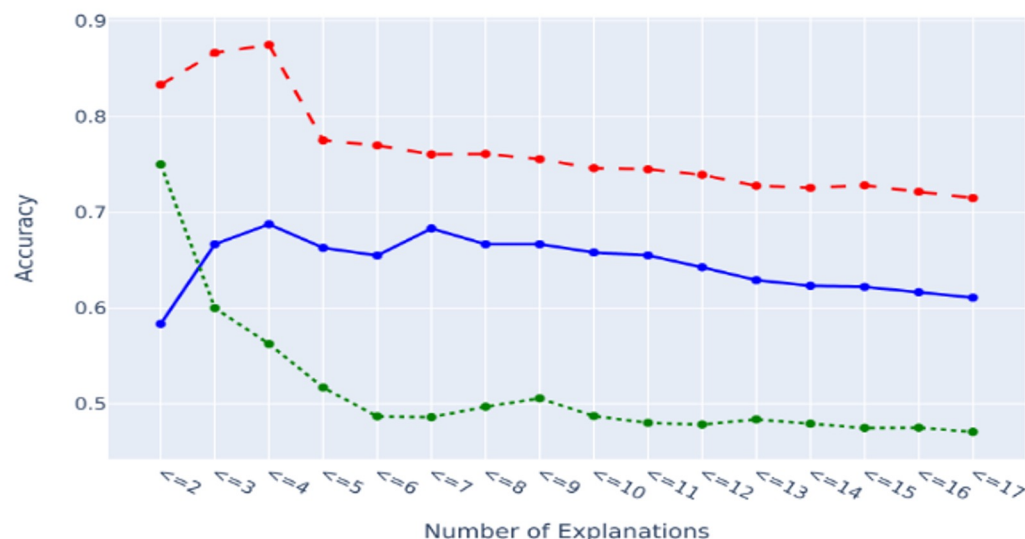$$C(f_j^{h_i}, f_k^{h_i}) = \frac{|t(f_j^{h_i}) \cap t(f_k^{h_i})|}{max(|t(f_j^{h_i})|, |t(f_k^{h_i})|)}$$

**Prior semantic/inference knowledge**

$$\omega_e(v_j, v_k; \theta_1) = \begin{cases} \theta_{gg}D(v_j, v_k) & v_j, v_k \in F_G^{h_i} \\ \theta_{aa}D(v_j, v_k) & v_j, v_k \in F_A^{h_i} \\ \theta_{ga}C(v_j, v_k) & v_j \in F_G^{h_i}, v_k \in F_A^{h_i} \\ \theta_{qg}C(v_j, v_k) & v_j \in F_G^{h_i}, v_k = h_i \\ \theta_{qa}C(v_j, v_k) & v_j \in F_A^{h_i}, v_k = h_i \end{cases}$$

$$\omega_v(v_i^{h_i}; \theta_2) = \begin{cases} \theta_{lr}L(v_j, h_i) + \theta_{ss}S(v_j, h_i) & v_j \in F_A^{h_i} \\ 0 & v_i \in F_G^{h_i} \\ 0 & v_i = h_i \end{cases}$$

Thayaparan, Valentino, Freitas, ACL Findings (2021)

| # | Approach | Accuracy | |
|---|----------|:--------:|:---:|
| | | WT | ARC |
| 1 | ExplanationLP (Best) | **61.37** | **40.21** |
| **Structure** | | | |
| 2 | Grounding-Abstract Categories | 58.33 | 35.13 |
| 3 | Edge weights | 43.78 | 29.45 |
| 4 | Node weights | 42.80 | 27.87 |
| **Cohesion** | | | |
| 5 | Hypothesis-Abstract cohesion | 38.71 | 30.37 |
| 6 | Hypothesis-Grounding cohesion | 59.33 | 38.72 |
| 7 | Grounding-Abstract cohesion | 59.12 | 38.14 |
| **Diversity** | | | |
| 8 | Abstract-Abstract diversity | 60.16 | 37.62 |
| 9 | Grounding-Grounding diversity | 60.44 | 37.71 |
| **Relevance** | | | |
| 10 | Hypothesis-Abstract semantic similarity | 55.38 | 35.49 |
| 11 | Hypothesis-Abstract lexical relevance | 54.68 | 36.01 |



red: ExplanationLP
blue: BERT$_{Large}$
green: PathNet

# of parameters:
- BERTBase: 110M parameters
- BERTLarge: 340M parameters
- ExplanationLP: 9 parameters

Thayaparan, Valentino, Freitas, ACL Findings (2021)

# Diff-Explainer: End-to-end abductive learning

An end-to-end differentiable framework that incorporates constraints via convex optimization layers into broader transformers-based architectures.

**Differentiable convex optimization (DCX) layers** (Agrawal et al., 2019) provide a way to encode constraints as part of a deep neural network.

**Problem:** ILP formulation is non-convex and cannot be incorporated into a differentiable convex optimization layer.

**Solution:**
- Approximate ILP with convex optimization constraints.
- Semi-Definite programming (SDP) is non-linear but convex and has shown to efficiently approximate combinatorial problems.
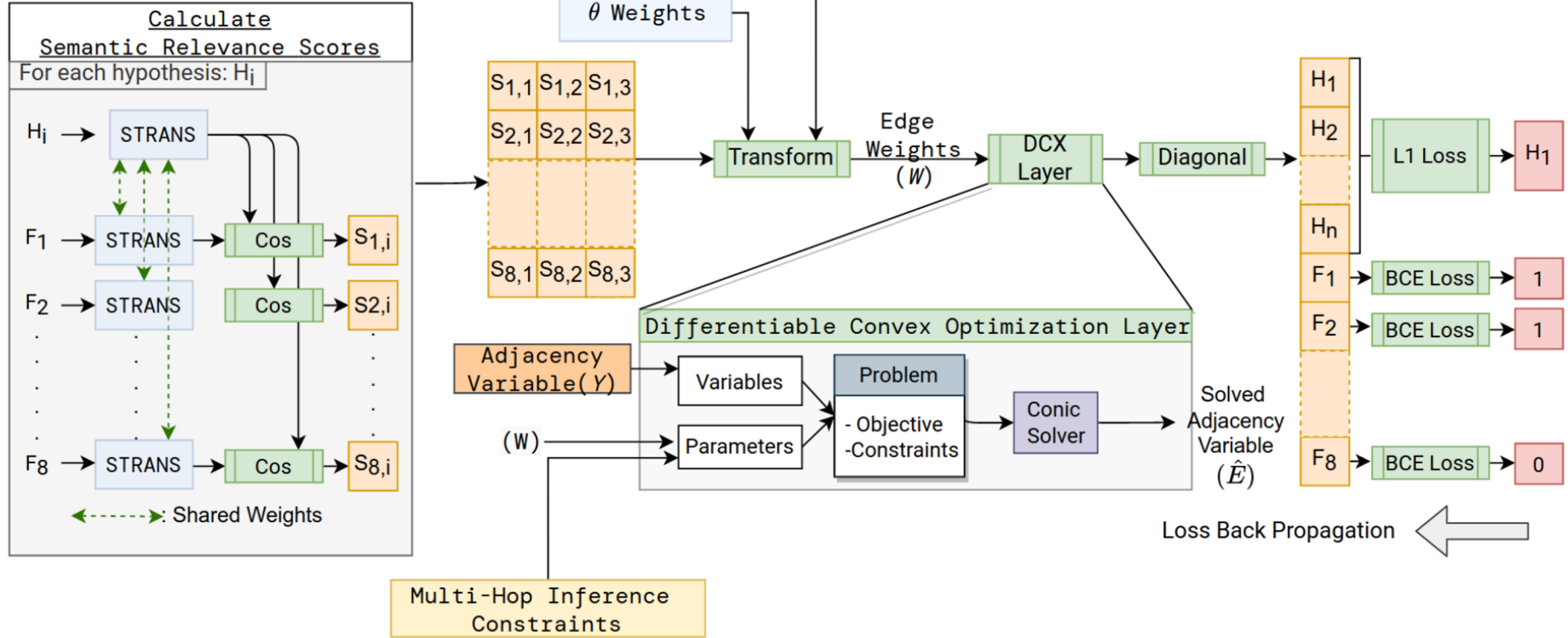
Specifically, we incorporate a differentiable convex optimization layer with Sentence-Transformers (STrans).

Thayaparan, Valentino, Ferreira, Rozanova, Freitas, TACL (2022)

# Programmable abductive NLI Solver

**Relevance**

$$s_{ij} = S(\vec{h_i}, \vec{f_j}) = \frac{\vec{h_i} \cdot \vec{f_j}}{\|\vec{h_i}\|\|\vec{f_j}\|}$$
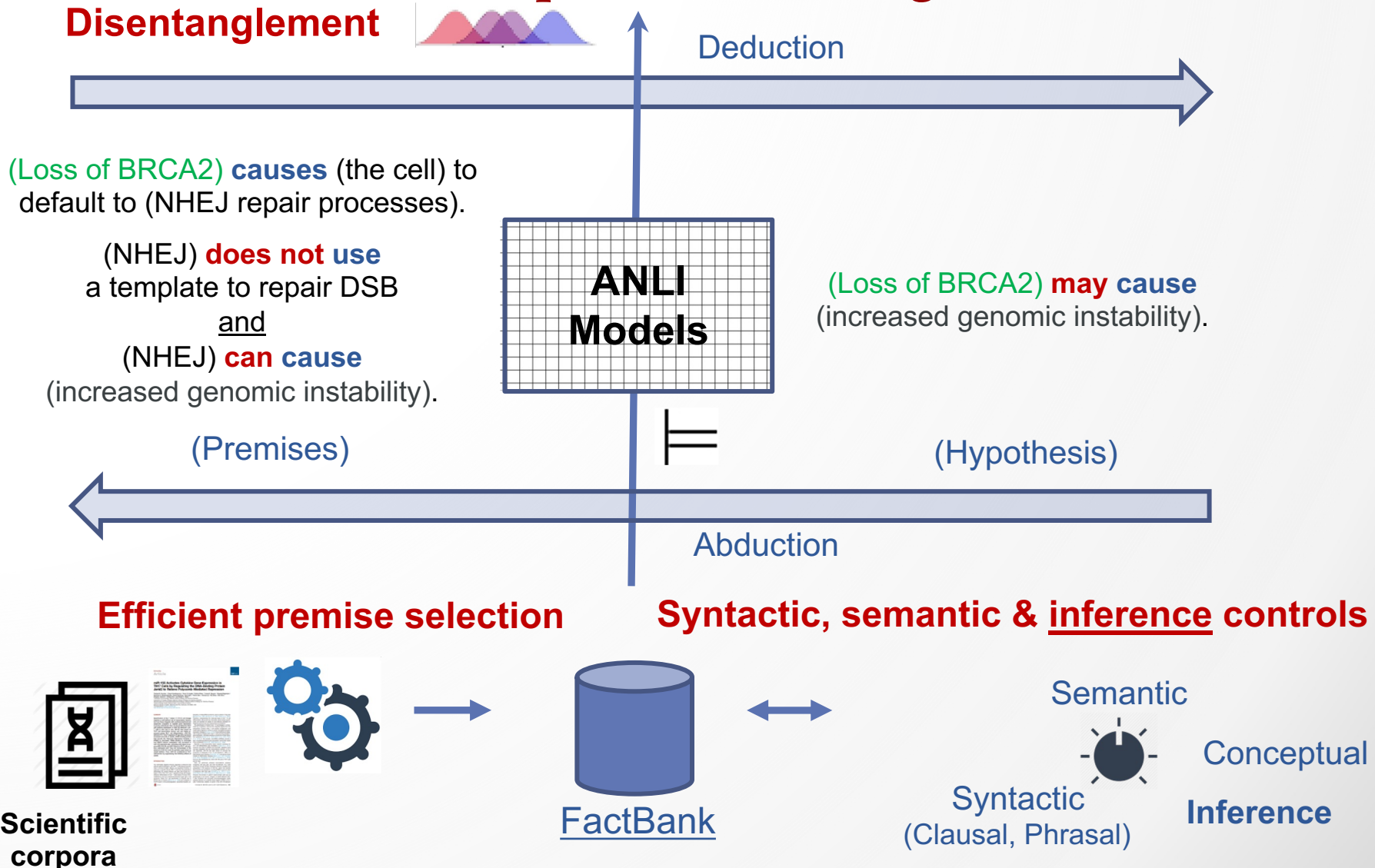
**Sentence embeddings**

Lexical Relevance Scores

$\theta$ Weights

### Calculate Semantic Relevance Scores

For each hypothesis: $H_i$

$H_i \rightarrow$ STRANS

$F_1 \rightarrow$ STRANS $\rightarrow$ Cos $\rightarrow$ $S_{1,i}$

$F_2 \rightarrow$ STRANS $\rightarrow$ Cos $\rightarrow$ $S_{2,i}$

$F_8 \rightarrow$ STRANS $\rightarrow$ Cos $\rightarrow$ $S_{8,i}$

⋯⋯➤ : Shared Weights

| $S_{1,1}$ | $S_{1,2}$ | $S_{1,3}$ |
|---|---|---|
| $S_{2,1}$ | $S_{2,2}$ | $S_{2,3}$ |
| $S_{8,1}$ | $S_{8,2}$ | $S_{8,3}$ |

Transform $\rightarrow$ Edge Weights $(W)$ $\rightarrow$ DCX Layer $\rightarrow$ Diagonal

### Differentiable Convex Optimization Layer

Adjacency Variable $(Y)$ $\rightarrow$ Variables

$(W) \rightarrow$ Parameters

Problem
- Objective
- Constraints

$\rightarrow$ Conic Solver $\rightarrow$ Solved Adjacency Variable $(\hat{E})$

$H_1$
$H_2$
$H_n$
$F_1$
$F_2$
$F_8$

$H_1 \rightarrow$ L1 Loss $\rightarrow$ $H_1$

$F_1 \rightarrow$ BCE Loss $\rightarrow$ 1

$F_2 \rightarrow$ BCE Loss $\rightarrow$ 1

$F_8 \rightarrow$ BCE Loss $\rightarrow$ 0

Loss Back Propagation

Multi-Hop Inference Constraints

**Saturation**

$$l_{ij} = L(h_i, f_j) = \frac{|trm(h_i) \cap trm(f_j)|}{max(|trm(h_i)|, |trm(f_j)|)}$$

**Prior semantic/inference knowledge**

$$W_{ij} = [\theta_1^s, \theta_2^s, \ldots, \theta_n^s] \cdot [s_{ij}^{\mathcal{D}_1}, s_{ij}^{\mathcal{D}_2}, \ldots, s_{ij}^{\mathcal{D}_n}]$$
$$+ [\theta_1^l, \theta_2^l, \ldots, \theta_n^l] \cdot [l_{ij}^{\mathcal{D}_1}, l_{ij}^{\mathcal{D}_2}, \ldots, l_{ij}^{\mathcal{D}_n}]$$

# Exploiting the structure of scientific explanations for

## multi-hop inference design

# Encoding abstract, mathematical inference

| Conjecture | Premise |
|---|---|
| Let $T = (S, \tau)$ be a topological space.<br>Let $A, B$ be subsets of $S$.<br>Then:<br>$\partial(A \cap B) \subseteq \partial A \cup \partial B$ where $\partial A$ denotes the boundary of $A$. | Let $S, T_1, T_2$ be sets such that $T_1, T_2$ are both subsets of $S$.<br>Then, using the notation of the relative complement:<br>$ST_1 \cap T_2 = ST_1 \cup ST_2$ |
| $\int \frac{X}{x(x^2-a^2)} = \frac{1}{2a^2}, \ln \frac{x^2-a^2}{x^2} + C$<br>for $x^2 > a^2$. | $\int \frac{dx}{x} = \ln x + C$<br>for $x \neq 0$. |
| Let $T = S, \tau$ be a compact space.<br>Then $T$ is countably compact. | Let $T = (S, \tau_{a,b})$ be a modified Fort space.<br>Then $T$ is not a $T_3$ space, $T_4$ space or $T_5$ space. |

*STAR: Cross-modal STAtement Representation for selecting relevant mathematical premises*

Ferreira & Freitas, EACL (2021)

*To be or not to be an Integer? Encoding Variables for Mathematical Text*

Ferreira et al., EACL (2021)

**Abstract statement representation**

*Similarity-based equational inference in physics*

Meadows & Freitas, PRR (2021)

*Premise Selection in Natural Language Mathematical Texts*

Ferreira & Freitas, ACL (2020)

**Multi-hop mathematical inference**

# Interventional, causal and granular evaluation of semantic and inference properties
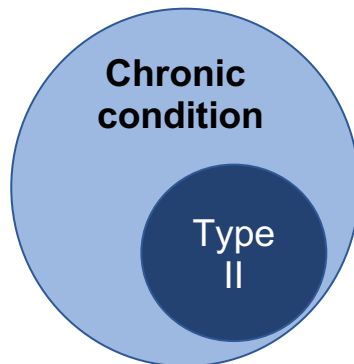


**ANLI Models**

Is reasoning really happening? (**quantifying causal effect**)
Are the semantic features present in the representations? (**probing**)
Do models reveal behavioural consistency? (**metamorphic testing**)

Inferences should follow **logical regularities** based on **abstract semantic features**.

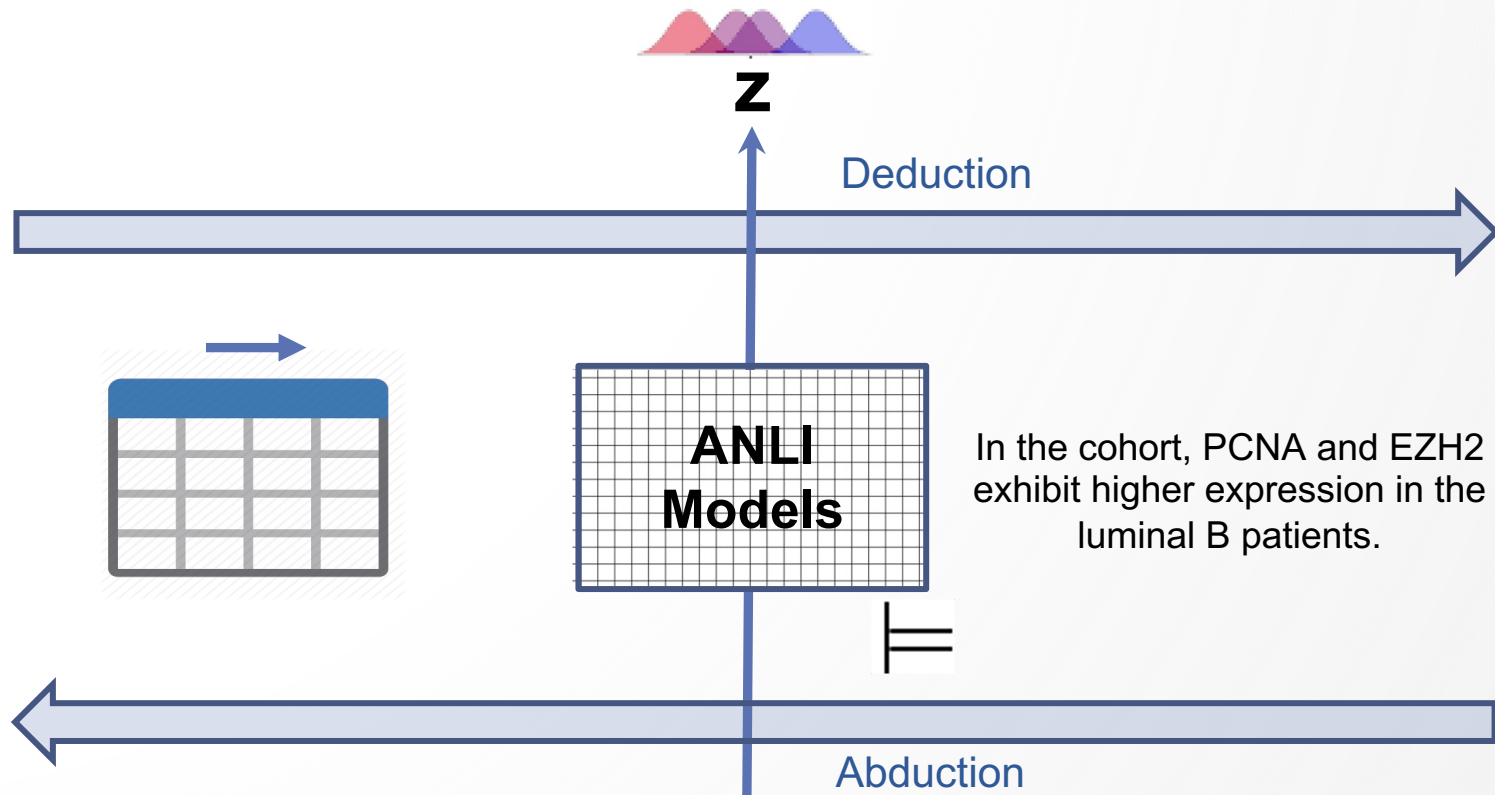"The patient does not have a **chronic condition**" $\models$ "The patient does not have **type II Diabetes** "



conceptual relations $\subseteq$ monotonicity

Manino et al. ACL Findings (2022)
Rozanova et al.  BlackBoxNLP@EMNLP (2022)
Rozanova et al., NALOMA (2021)

# Encoding inferences over table evidence



Badaro, Saeed, Papotti, TACL (2023)

# Take away

Building natural language explanation machines for science (key neuro-symbolic strategies)

- Granular, controlled, neuro-symbolic inference.
- Integrating and controlling LLMs properties.
- **Sentence-level encoding strategies:**
  - Exploit disentanglement with semantic priors.
  - Use robust lightweight semantic representations.
  - Organise your semantic space.
  - Use specialised strategies for concepts and abstract statements.
- **Inference-level encoding strategies:**
  - Exploit the structure of scientific explanations.
  - Integrate inference priors as constraints.
  - Use differentiable convex layers for end-to-end training.
- **Evaluation of true inference properties:**
  - Determine true reasoning performance by using causal, interventional methods.
  - Probe your model for key properties.
  - Apply systematic behavioural testing.

Thank you

The northern hemisphere **is a kind** of hemisphere of earth

Abstraction (hemisphere)

a hemisphere of earth i**s a kind** of place

Abstraction (place)

If a place **is in** summer, then it **will have** the most sunlight

Unification

Northern hemisphere **will have** the most sunlight in summer

# Granular evaluation



Ferreira et al., ACL Demo (2021)

*Does My Representation Capture X?*
*Probe-Ably*

**Question:** A group of students are studying bean plants. All of the following traits are affected by changes in the environment except . . .

**Candidate answers:** [A] leaf color [B] seed type [C] bean production [D] plant height

**Explanation**

(i) The type of seed of a plant is an inherited characteristic;

(ii) Inherited characteristics are the opposite of acquired characteristics;

(iii) An organism's environment affects that organism's acquired characteristics;

(iv) A plant is a kind of organism;

(v) A bean plant is a kind of plant;

(vi) Trait is synonymous with characteristic.

(i) The type of seed of a plant is an inherited characteristic;

$\forall xy(plant(x) \land seedType(y, x) \rightarrow characteristic(y, x) \land inherited(y))$

(ii) Inherited characteristics are the opposite of acquired characteristics;

$\forall xy(characteristic(x, y) \land inherited(x) \rightarrow \neg acquired(x))$

(iii) An organism's environment affects that organism's acquired characteristics;

$\forall xyw(organism(x) \land environment(y, x) \land characteristic(w, x) \land acquired(w) \rightarrow$

$\rightarrow \exists e(affect(e) \land agent(e, y) \land patient(e, w))$

(iv) A plant is a kind of organism;

$\forall x(plant(x) \rightarrow organism(x))$

(v) A bean plant is a kind of plant;

$\forall x(beanPlant(x) \rightarrow plant(x))$

(vi) Trait is synonymous with characteristic.

$\forall xy(trait(x, y) \leftrightarrow characteristic(x, y))$

**KB Φ**

$\forall xy(plant(x) \land seedType(y, x) \rightarrow characteristic(y, x) \land inherited(y))$

$\forall xy(characteristic(x, y) \land inherited(x) \rightarrow \neg acquired(x))$

$\forall xyw(organism(x) \land environment(y, x) \land characteristic(w, x) \land acquired(w) \rightarrow$
$\rightarrow \exists e(affect(e) \land agent(e, y) \land patient(e, w))$

$\forall x(plant(x) \rightarrow organism(x))$

$\forall x(beanPlant(x) \rightarrow plant(x))$

$\forall xy(trait(x, y) \leftrightarrow characteristic(x, y))$

**Φ ⊨ ψ ?**

**Question:** find a characteristic of plants not affected by those plants' environments. That is, we are asked for a **P** making the schematic formula **true**.

**ψ**

$\forall xyzwe(beanPlant(x) \land environment(y, x) \land changeIn(z, y) \land trait(w, x) \land affect(e) \land agent(e, z) \land \mathbf{P} \rightarrow \neg patient(e, w))$

**P**: $seedType(w, x)$

Valentino, Pratt-Hartmann, Freitas, IWCS (2021)

# Question: How a 1 degree rise in temperature will affect the grape harvest in Valais?

Evidence-based explanation

| Optimal Wine Grape Temperatures Over the Growing Season | | |
|---|---|---|
| **Variety** | **Min** | **Max** |
| Pinot Gris | 13°C | 15°C |
| Riesling | 13°C | 17°C |
| Pinot Noir | 14°C | 15°C |
| Chardonnay | 14°C | 18°C |
| Sauvignon Blanc | 14°C | 18°C |
| Syrah | 16°C | 19°C |
| Table Grapes | 19°C | 22°C |

Fine wine production is likely to shift due to climate change. Among agricultural products, wine grapes are one of the most sensitive crops to variations in temperature and precipitation

Since the year 1864, the temperature in the Canton of Valais has increased by 2 °C. If global greenhouse gas emissions continue to rise in the future, the warming will continue and will amount to further 3 °C by 2060 with respect to the mean of the period 1981-2010.

Pinot Noir accounts for 11% of the grape production in Valais.

| Availability by Grape/Blend =-Valais | |
|---|---|
| Pinot Noir | 11% |
| Petite Arvine | 9% |
| Chasselas | 9% |
| Syrah | 8% |
| Rare Red Blend | 7% |
| Corlanin | 6% |
| Gamay Pinot Noir | 6% |
| Other | 43% |

| Temperature Deviation - Valais (°C) | | RCP2.6 | | | RCP8.5 | | |
|---|---|---|---|---|---|---|---|
| | | min | mid | max | min | mid | max |
| Summer | 2035 | 0.8 | 1.8 | 2.5 | 1.7 | 2.0 | 3.0 |
| | 2060 | 1.0 | 1.9 | 3.1 | 2.7 | 3.9 | 5.7 |
| Winter | 2035 | 0.6 | 1.0 | 1.8 | 0.9 | 1.7 | 1.9 |
| | 2060 | 0.8 | 1.5 | 1.9 | 1.8 | 2.2 | 2.9 |

| Precipitation Deviation - Valais (%) | | RCP2.6 | | | RCP8.5 | | |
|---|---|---|---|---|---|---|---|
| | | min | mid | max | min | mid | max |
| Summer | 2035 | -18 | -1 | 5 | -10 | -3 | 4 |
| | 2060 | -15 | 0 | 12 | -20 | -12 | 12 |
| Winter | 2035 | -10 | 8 | 20 | 2 | 15 | 22 |
| | 2060 | -3 | 12 | 20 | 0 | 12 | 22 |

**Claim:** US Summer Youth Employment programs can be replicated in Central Europe to help low income youth overcome barriers to accessing jobs.

**Labor**

**Criminal Justice System**

**Education**

SYEPs disproportionately serve youth from low-income households that typically face higher than average barriers to entering the labor market

SYEPs consistently reduce involvement in the criminal justice system for participating youth for the duration of the program and at least a year beyond.

Evidence on SYEPs' role in improving educational outcomes is mixed.

For the most part, SYEPs do not increase rates of formal sector employment for the average participant after the program ends, with some exceptions

Youth at greater risk of experiencing socially costly outcomes, such as involvement with the criminal justice system or disengagement from school, are shown to experience the greatest benefits from SYEP.

On average, in the studies that showed positive effects in academic outcomes, those who benefited were youth of legal drop-out age and youth who had a higher rate of school absences prior to program participation.

|  | New York City 2006 – 2010 | Boston 2015 | Chicago 2012 |
|---|---|---|---|
| Youth offered SYEP job | 72.3% | 83.6% | 78.7% |
| Youth not offered SYEP job | 18.5% | 26.4% | 15.2% |

The move toward evidence-based policy (EBP) formation still requires improvement of the understanding of the role of evidence within policy process and analysis of the barriers in using evidence in policy development processes

SYEP participation decreases arrests and convictions during the program summer

| | Control | Participants | Impact |
|---|---|---|---|
| Violent crime arrests per hundred youth | 18.34% | 11.96% | -6.38% |

Outcomes suggests improvements in social-emotional skills, academic and career aspirations, and work habits associated with job readiness.

**Claim:** it is not necessary to prove the absence of the debtor's assets to obtain the disregard of legal personality.

**Precedent Nº 1.729.554 (Superior Court)**: "In fact, the disregard of the legal personality can be decreed even if insolvency is not configured, provided that the deviation of purpose or the patrimonial confusion, characterizing the abuse of personality, are verified."

**Law 13105 (Federal - CPC):  Art. 134**: The incident of disregard is applicable at all stages of the acknowledgement process, in the execution of the sentence and in the execution based on an extrajudicial enforcement order

**Law 13105 (Federal - CPC) Art. 134 § 4:** The application must demonstrate the completion of the specific legal presuppositions for disregarding the legal personality.

**Law 10406/02 (Federal – Civil code) Art. 50:** In case of abuse of legal personality, characterized by the misuse of purpose, or by the confusion of assets, the judge may decide, at the request of the party, or of the Public Prosecutor's Office when it is up to him to intervene in the process, that the effects of certain and certain relationships of obligations are extended to the private assets of the administrators or partners of the legal entity.

**Doctrine (Humberto Dalla)**: "In the case of greater disregard, in which the true passive holder of the credit is the partner (who acted abusively through the legal entity), the author has the right to choose his liability, regardless of the potential satisfaction of the credit before the legal entity."

# Complex Sentence Representation

- Clausal-Phrasal Disembedding (CPD).
- Minimal, localised, self-contained propositions.

"Programmed death-ligand 1 (PD-L1) also known as cluster of differentiation 274 (CD274) or B7 homolog 1 (B7-H1) is a protein that in humans is encoded by the CD274 gene."

Core: PD-L1 is encoded by the CD274 gene.
Context: This is in humans.

PD-L1 is also known as cluster of differentiation 274.
PD-L1 is also known as B7-H1.
PD-L1 is a protein.

Programmed death-ligand 1 has abbreviation PD-L1.
Cluster of differentiation 274 has abbreviation CD274.
B7 homolog 1 has abbreviation B7-H1.

Niklaus, Cetto, Freitas, Handschuh ACL (2019)

# Complex Sentence Representation

"Programmed death-ligand 1 (PD-L1) also known as cluster of differentiation 274 (CD274) or B7 homolog 1 (B7-H1) is a protein that in humans is encoded by the CD274 gene."

Core: is encoded by(e, PD-L1, the CD274 gene).
Context: in(e, humans).

is also known as(PD-L1, cluster of differentiation 274).
is also known as(PD-L1, B7-H1).
is a(PD-L1, protein).

has abbreviation(Programmed death-ligand 1, PD-L1).
has abbreviation(Cluster of differentiation 274, CD274).
has abbreviation(B7 homolog 1, B7-H1).

# Complex Sentence Representation

| | CLAUSAL/PHRASAL TYPE | HIERARCHY | # RULES |
|---|---|---|---|
| | **Clausal disembedding** | | |
| 1 | Coordinate clauses | coordinate | 1 |
| 2 | Adverbial clauses | subordinate | 6 |
| 3a | Relative clauses (non-restrictive) | subordinate | 5 |
| 3b | Relative clauses (restrictive) | subordinate | 4 |
| 4 | Reported speech | subordinate | 4 |
| | **Phrasal disembedding** | | |
| 5 | Coordinate verb phrases | coordinate | 1 |
| 6 | Coordinate noun phrases | coordinate | 2 |
| 6 | Participial phrases | subordinate | 4 |
| 8a | Appositions (non-restrictive) | subordinate | 1 |
| 8b | Appositions (restrictive) | subordinate | 1 |
| 9 | Prepositional phrases | subordinate | 3 |
| 10 | Adjectival and adverbial phrases | subordinate | 2 |
| 11 | Lead NPs | subordinate | 1 |
| | Total | | 35 |

| System | Precision | Recall | $F_1$ | AUC |
|---|---|---|---|---|
| REVERB | -7.2% | +19.1% | +8.4% | **+39.2%** |
| OLLIE | +1.1% | -1.5% | -0.3% | -1.1% |
| ClausIE | +17.0% | -3.5% | +8.1% | +13.0% |
| Stanford Open IE | **+25.0%** | **+27.2%** | **+25.5%** | +35.5% |
| PropS | -6.1% | +16.9% | +4.5% | +12.4% |
| OpenIE-4 | +10.0% | +8.6% | +9.4% | +19.6% |
| MinIE | +23.7% | -1.8% | +12.5% | +21.3% |
| OpenIE-5 | +5.0% | +4.2% | +4.6% | +9.0% |
| RnnOIE | -15.0% | +0.9% | -8.3% | -14.1% |

Niklaus, Cetto, Freitas, Handschuh ACL (2019)