# Exploring structures of inferential mechanisms through simplistic digital circuits

25 October 2025, AIC workshop @ ECAI 2025

*10th International Workshop on Artificial Intelligence and Cognition*

Giovanni Sileno
g.sileno@uva.nl

*University of Amsterdam*

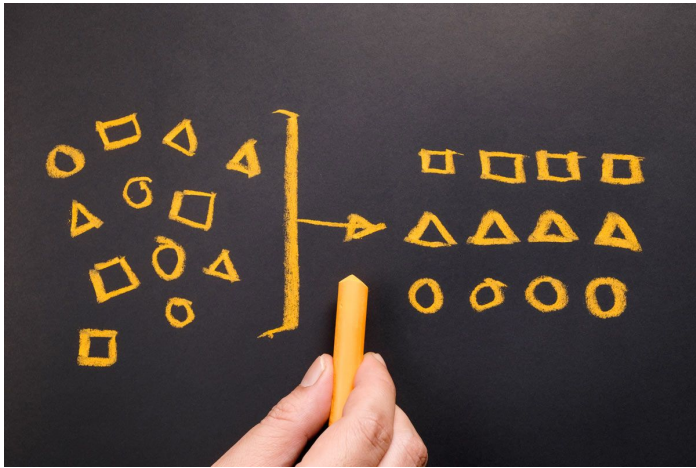Jean-Louis Dessalles
jean-louis.dessalles@telecom-paris.fr

*Télécom Paris*

# Inferential mechanisms in cognitive studies

- Cognitive studies have since long distinguished several types of cognitive mechanisms, by conducting distinct modelling efforts and different types of experiments.

# Inferential mechanisms in cognitive studies

- Cognitive studies have since long distinguished several types of cognitive mechanisms, by conducting distinct modelling efforts and different types of experiments.



**categorization**

the process grouping objects, events, or situations, on the basis of shared characteristics

# Inferential mechanisms in cognitive studies

● Cognitive studies have since long distinguished several types of cognitive mechanisms, by conducting distinct modelling efforts and different types of experiments.



**categorization**

the process grouping objects, events, or situations, on the basis of shared characteristics

*via prototype theory, exemplar theory, rule-based, knowledge-based, Bayesian models…*

# Inferential mechanisms in cognitive studies

- Cognitive studies have since long distinguished several types of cognitive mechanisms, by conducting distinct modelling efforts and different types of experiments.



**induction**

the process drawing general
rules/models from observations

# Inferential mechanisms in cognitive studies

- Cognitive studies have since long distinguished several types of cognitive mechanisms, by conducting distinct modelling efforts and different types of experiments.



**induction**

the process drawing general
rules/models from observations

*via associative (co-occurrence), descriptive (similarity), Bayesian models, …*

# Inferential mechanisms in cognitive studies

- Cognitive studies have since long distinguished several types of cognitive mechanisms, by conducting distinct modelling efforts and different types of experiments.

**deductive reasoning**

**categorization**
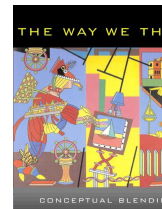
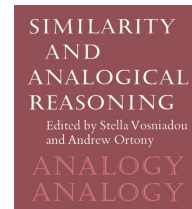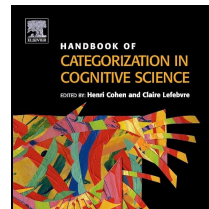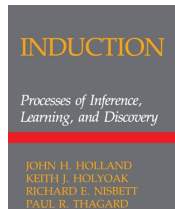**causal/diagnostic inference**

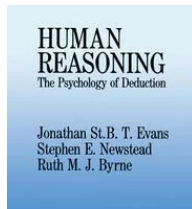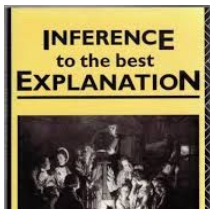**analogical reasoning**

**abduction**
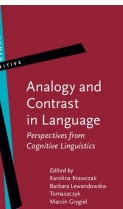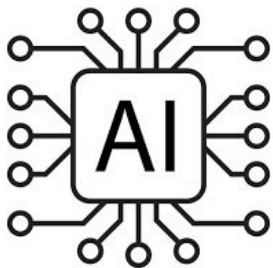
**contrast**

**induction**

**conceptual merge/blending**

….

# Inferential mechanisms in AI

- The same dispersion can be observed TRADITIONALLY in artificial intelligence, research and practice….
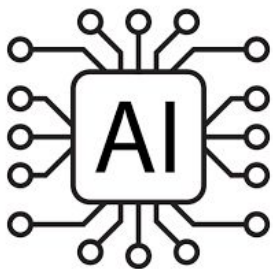


**categorization**

the process grouping objects, events, or situations, on the basis of shared characteristics

**induction**

the process drawing general rules/models from observations

# Inferential mechanisms in AI

- The same dispersion can be observed TRADITIONALLY in artificial intelligence, research and practice….

**categorization**

the process grouping objects, events, or situations, on the basis of shared characteristics

**symbolic AI:** *rule-based systems, decision trees, formal concept analysis (FCA), …* **sub-symbolic AI:** *neural networks, support vector machines, clustering, …*

**induction**

the process drawing general rules/models from observations

**symbolic AI:** *inductive logic programming, version-space, and explanation-based learning*
**sub-symbolic IA:** *all machine learning methods (including deep learning and generative AI methods)*

# Inferential mechanisms in AI

- The same dispersion can be observed TRADITIONALLY in artificial intelligence, research and practice…. TODAY, with transformer architectures, there is a general belief that we can induce essentially all inferential mechanisms having adequate data, so the problem is now to extract them!

# Inferential mechanisms in AI

- The same dispersion can be observed TRADITIONALLY in artificial intelligence, research and practice…. TODAY, with transformer architectures, there is a general belief that we can induce essentially all inferential mechanisms having adequate data, so the problem is now to extract them!

**MECHANISTIC INTERPRETABILITY**

*reverse engineer* the neural network to identify "circuits" implementing inferential functions

Observed model

Hypothetical disentangled model

Bereska, Leonard, and Efstratios Gavves. "Mechanistic interpretability for AI safety--a review." *arXiv preprint arXiv:2404.14082* (2024). https://leonardbereska.github.io/blog/2024/mechinterpreview/
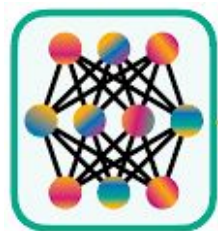
# Inferential mechanisms in AI

- The same dispersion can be observed TRADITIONALLY in artificial intelligence, research and practice…. TODAY, with transformer architectures, there is a general belief that we can induce essentially all inferential mechanisms from adequate data, so the problem is now to extract them!
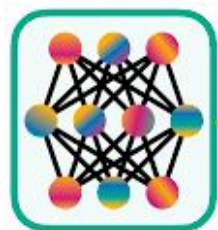
**MECHANISTIC INTERPRETABILITY**

*reverse engineer* the neural network to identify "circuits" implementing inferential functions



Observed model → Hypothetical disentangled model

affinity with Computational Neuro-Science

# Inferential mechanisms in AI

- The same dispersion can be observed TRADITIONALLY in artificial intelligence, research and practice…. TODAY, with transformer architectures, there is a general belief that we can induce essentially all inferential mechanisms from adequate data, so the problem is now to extract them!
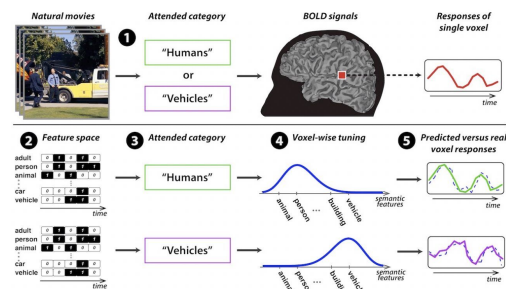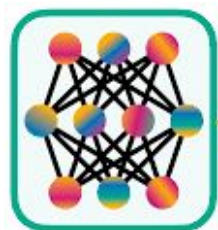


**MECHANISTIC INTERPRETABILITY**

*reverse engineer* the neural network to identify "circuits" implementing inferential functions

affinity with Computational Neuro-Science

**still none of these approaches say where all inferential mechanisms come from!**

# Restarting from scratch…

- Rather then reverse engineering the brain, we could investigate an inferential system from its minimal core. But what would that be?

# Restarting from scratch…

- Rather then reverse engineering the brain, we could investigate an inferential system from its minimal core. But what would that be?

- **Intuition:** several inferential mechanisms have been reproduced with success through symbolic AI methods ⇒ they *SHOULD* get some aspect of the cognitive functions right.

# Restarting from scratch…

- Rather then reverse engineering the brain, we could investigate an inferential system from its minimal core. But what would that be?

- **Intuition:** several inferential mechanisms have been reproduced with success through symbolic AI methods ⇒ they *SHOULD* get some aspect of the cognitive functions right.

- **This is a weaker assumption than hypotheses like *Language of Thought* (LoT) or *Physical Symbol System* (PSS)!**

# Outline of presentation

- Under which constraints logic rules can be interpreted as valid digital circuits?

  - Which inferential constructs can be constructed from these constraints?

    - Generalizing this analysis, we will find an unexpectedly unifying schema!

# Let's take a classic symbolic tool for computational inference: logic programming

```
…
parent(marge, lisa).
parent(marge, bart).
parent(marge, maggie).
parent(homer, lisa).
parent(homer, bart).
parent(homer, maggie).
parent(abraham, homer).
parent(abraham, herb).
parent(mona, homer).


child(X,Y) :- parent(Y,X).

?- child(lisa, marge).
true
```

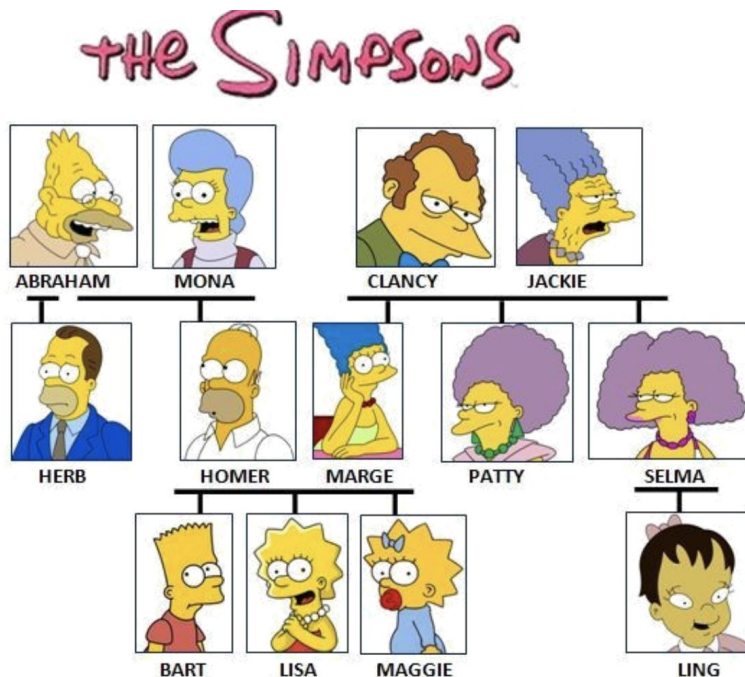# Let's take a classic symbolic tool for computational inference: logic programming

```
…
parent(marge, lisa).
parent(marge, bart).
parent(marge, maggie).
parent(homer, lisa).
parent(homer, bart).
parent(homer, maggie).
parent(abraham, homer).
parent(abraham, herb).
parent(mona, homer).


orphan(X) :- not parent(Y, X).


?- orphan(abraham).
true
```
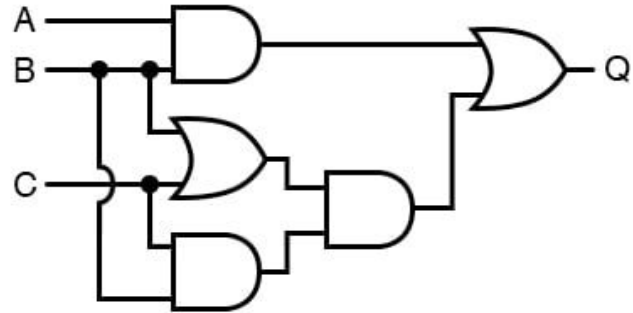
**DEFAULT NEGATION**
creating
knowledge
out of "ignorance"

# Going electrical!

- In order to drop implicit assumptions holding with symbolic methods, let's think in terms of **digital circuits**: activation at the end maps just to an electrical feedforward mechanism!

# Going electrical!

- It seems easy to associate logic rules to logic ports.

```
p :- a, b.
```

# Going electrical!

- It seems easy to associate logic rules to logic ports.

```
p :- a, b.
```



a
b
p
**AND**

- Yet, in propositional logic, the **contrapositive** holds:

$$a \wedge b \to p \quad \Longleftrightarrow \quad \neg p \to \neg a \vee \neg b$$

**CLASSIC NEGATION**
assertions of denial (*it is the case that not …*)

# Going electrical!

- It seems easy to associate logic rules to logic ports.

```
p :- a, b.
```


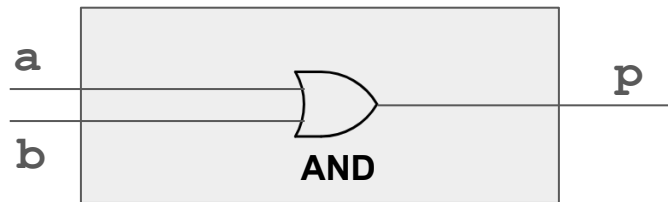
a

b

**AND**

p

- Yet, in propositional logic, the **contrapositive** holds:

$$a \wedge b \rightarrow p \qquad \Longleftrightarrow \qquad \neg p \rightarrow \neg a \vee \neg b$$

**CLASSIC NEGATION**

assertions of denial (*it is the case that not …*)

**This dual circuit is not directly implementable!**
The output is not a single port, and it is non-deterministic.

# Going electrical!

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.
```



$a \wedge b \rightarrow p$ ⬌ $\neg p \rightarrow \neg a \vee \neg b$

The only valid circuits can be constructed
**removing this indeterminism**.

# Going electrical!

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.
```

$a \wedge b \rightarrow p$ ⬅➡ $\neg p \rightarrow \neg a \vee \neg b$

The only valid circuits can be constructed **removing this indeterminism**.

# Going electrical!

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.
```

$a \wedge b \rightarrow p$ ⬌ $\neg p \rightarrow \neg a \vee \neg b$

Semantically, however, logic interacts with the space of **all possible models/states.**

# Going electrical!

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.
```

$a \wedge b \rightarrow p$ ⟺ $\neg p \rightarrow \neg a \vee \neg b$

Semantically, however, logic interacts with the space of **all possible models/states.**

*How to recreate this electrically?*

# Going electrical!

deterministic machinery
reifying constraints

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```

combinatorial exploration

# What we learn by going electrical? (1)

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```



expressions do not have memory!

# What we learn by going electrical? (1)

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```



expressions do not have memory!

a

b

b

¬b

a

¬a

p

p

¬p

AND

AND

AND

memory/
state is
on the
channels

# What we learn by going electrical? (2)

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```



side-effects are related to the **data-flow**

memory/state is on the channels

# What we learn by going electrical? (2)



side-effects are
related to the
**data-flow**

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```

**conditionals are not simple operators!**
**implications are topological!**

a

b

b

¬b

a

¬a

p

p

¬p

AND

AND

AND

memory/
state is
on the
channels

# What we learn by going electrical? (3)

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```



channels may receive
opposite signals: *invalid status*

memory/
state is
on the
channels

# What we learn by going electrical? (3)

```
p :- a, b.
-a :- -p, b.
-b :- -p, a.

1{a; -a}1.
1{b; -b}1.
1{p; -p}1.
```

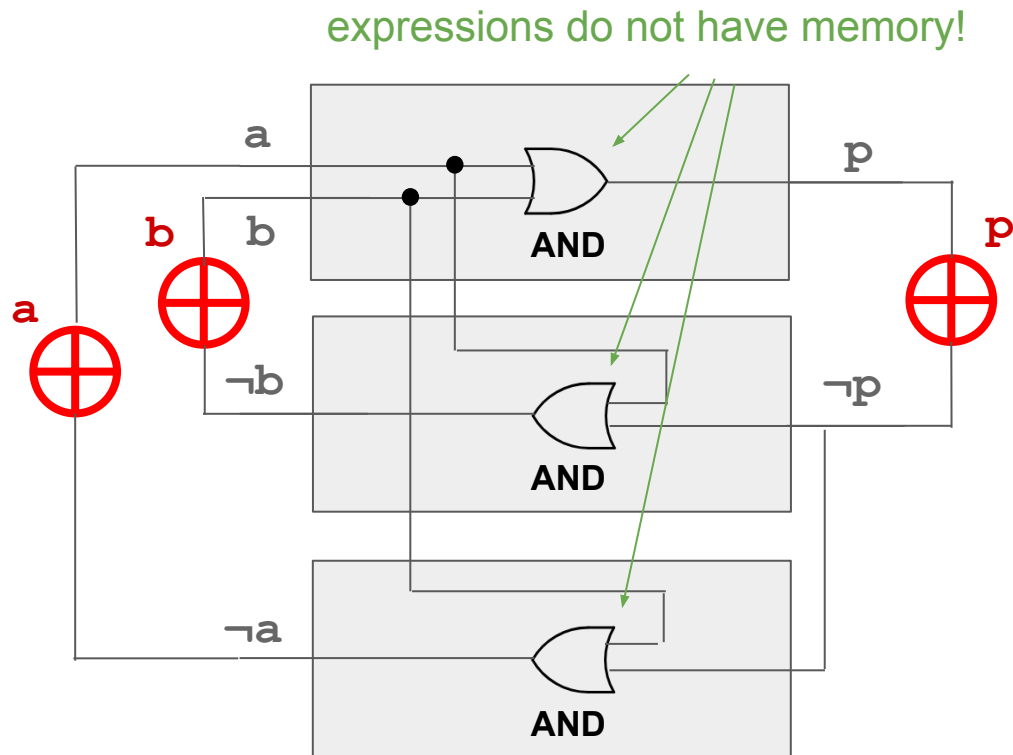**we should purge all models with contradictions**

channels may receive opposite signals: *invalid status*

a

b    b

a

¬b

¬a

AND

AND

AND

p

p

¬p

memory/ state is on the channels

# Going electrical: conclusions!

```
p :- a, b.
```



➡️ **logic rules** as in LP are more sound specifications of activation mechanisms. Yet, we should not use default negation!

# Going electrical: conclusions!

```
p :- a, b.
```



→ **logic rules** as in LP are more sound specifications of activation mechanisms. Yet, we should not use default negation!

→ We have some indication on **minimal activation mechanisms: inputs can be composed only with OR ";" and AND "," ..** *but what about outputs?*

# A suite of dependencies! (abusing the standard LP notation)

```
p :- a, b.
p :- a; b.
p, q :- a.
p; q :- a.
```

# 1. Conjunction in body

```
p :- a, b.
```

# 2. Disjunction in body

```
p :- a; b.
```

⟷

```
p :- a.
p :- b.
```

this acts like an
accumulation bus

a    $P_a$

b    $P_b$

p

# 3. Conjunction in head

p, q :- a. $\iff$ p :- a.
q :- a.

this acts like a broadcast bus

This form is also expressed in the contrapositive of (2)

a ∧ b → p $\iff$ ¬p → ¬a ∨ ¬b

# 4. Disjunction in head

`p; q :- a.`

This form is also expressed in the
contrapositive of (1)

$a \lor b \to p$ ⟺ $\neg p \to \neg a \land \neg b$

*non-deterministic!*



a → p
   → q

(If made deterministic, it would
be just as form 3.)

# 4. Disjunction in head

```
p; q :- a.
```

This form is also expressed in the contrapositive of (1)

a ∨ b → p ⟺ ¬p → ¬a ∧ ¬b

*non-deterministic!*



(If made deterministic, it would be just as form 3.)

**But not if interpreted as a XOR!**

# From propositions to predicates…

- So far we considered only the propositional case. However, features are always about some entity, just as predicates are always about something.

- Let us consider how these dependencies would appear in common usages in logic programming.

- We will consider two cases:

  - **Unary predicates** (concerning only a single entity)
  - **Binary predicates** (relating two distinct entities)

# Unary predicates (abusing the standard LP notation)

```
p(X) :- a(X), b(X).
p(X) :- a(X); b(X).
p(X), q(X) :- a(X).
p(X); q(X) :- a(X).
```

# Unary predicates

```
p(X) :- a(X), b(X).
p(X) :- a(X); b(X).
p(X), q(X) :- a(X).
p(X); q(X) :- a(X).
```



(1) relation defining new concepts:
```
angrydog(X) :- dog(X), angry(X).
```
(2) relation defining taxonomical relations:
```
mammal(X) :- dog(X); cat(X).
```
(3) relation activating back the source concepts:
```
dog(X), angry(X) :- angrydog(X).
```
(4) non-deterministic relation activating alternative choices:
```
dog(X); cat(X) :- mammal(X).
```

# Unary predicates

```
p(X) :- a(X), b(X).
p(X) :- a(X); b(X).
p(X), q(X) :- a(X).
p(X); q(X) :- a(X).
```

(1) relation defining new concepts:
    angrydog(X) :- dog(X), angry(X).

(2) relation defining taxonomical relations:
    mammal(X) :- dog(X); cat(X).

(3) relation activating back the source concepts:
    dog(X), angry(X) :- angrydog(X).

(4) non-deterministic relation activating alternative choices:
    dog(X); cat(X) :- mammal(X).


MAMMALS

# Unary predicates



```
p(X) :- a(X), b(X).
p(X) :- a(X); b(X).
p(X), q(X) :- a(X).
p(X); q(X) :- a(X).
```

(1) relation defining new concepts:
```
angrydog(X) :- dog(X), angry(X).
```
(2) relation defining taxonomical relations:
```
mammal(X) :- dog(X); cat(X).
```
(3) relation activating back the source concepts:
```
dog(X), angry(X) :- angrydog(X).
```
(4) non-deterministic relation activating alternative choices:
```
dog(X); cat(X) :- mammal(X).
```

# Unary predicates


MAMMALS

```
p(X) :- a(X), b(X).
p(X) :- a(X); b(X).
p(X), q(X) :- a(X).
p(X); q(X) :- a(X).
```

(1) relation defining new concepts:
```
    angrydog(X) :- dog(X), angry(X).
```
(2) relation defining taxonomical relations:
```
    mammal(X) :- dog(X); cat(X).
```
(3) relation activating back the source concepts:
```
    dog(X), angry(X) :- angrydog(X).
```
(4) non-deterministic relation activating alternative choices:
```
    dog(X); cat(X) :- mammal(X).
```

# Binary predicates



- Give an entity, we may bind it to other entities

$$dog(x) \wedge tail(y) \wedge has(x, y)$$

# Binary predicates



- Give an entity, we may bind it to other entities

$$dog(x) \land tail(y) \land has(x, y)$$

- At rule level, this becomes an *existential rule*, which is **not treatable** by standard logic programming derivation (nor by description logic reasoners)

$$\forall x : dog(x) \rightarrow \exists y : tail(y) \land has(x, y)$$

# Binary predicates



- Give an entity, we may bind it to other entities

$$dog(x) \land tail(y) \land has(x, y)$$

- At rule level, this becomes an *existential rule*, which is **not treatable** by standard logic programming derivation (nor by description logic reasoners)

$$\forall x : dog(x) \rightarrow \exists y : tail(y) \land has(x, y)$$



*there is a similar problem with active rules (causal laws!)*

# Binary predicates (abusing the standard LP notation)

```
p(X) :- Y/ a(X, Y), Z/ b(X, Z).
p(X) :- Y/ a(X, Y); Z/ b(X, Z).
Y/ p(X, Y), Z/ q(X, Z) :- a(X).
Y/ p(X, Y); Z/ q(X, Z) :- a(X).
```

# Binary predicates



```
p(X) :- Y/ a(X, Y), Z/ b(X, Z).
p(X) :- Y/ a(X, Y); Z/ b(X, Z).
Y/ p(X, Y), Z/ q(X, Z) :- a(X).
Y/ p(X, Y); Z/ q(X, Z) :- a(X).
```

(1) relation determining a concept by composition:
```
    car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
    student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
    Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
    Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

# Binary predicates



```
p(X) :- Y/ a(X, Y), Z/ b(X, Z).
p(X) :- Y/ a(X, Y); Z/ b(X, Z).
Y/ p(X, Y), Z/ q(X, Z) :- a(X).
Y/ p(X, Y); Z/ q(X, Z) :- a(X).
```



(1) relation determining a concept by composition:
```
    car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
    student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
    Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
    Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

# Binary predicates



```
p(X) :- Y/ a(X, Y), Z/ b(X, Z).
p(X) :- Y/ a(X, Y); Z/ b(X, Z).
Y/ p(X, Y), Z/ q(X, Z) :- a(X).
Y/ p(X, Y); Z/ q(X, Z) :- a(X).
```

(1) relation determining a concept by composition:
```
    car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
    student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
(3) reifies mereonomical relations (similarly to a **constructor in Object Oriented Programming**):
```
    Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
    Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```
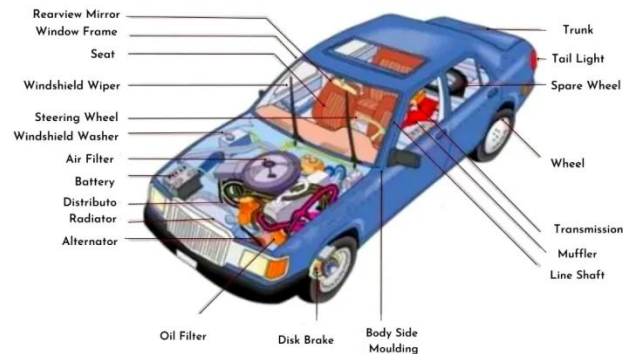
# Binary predicates





```
p(X) :- Y/ a(X, Y), Z/ b(X, Z).
p(X) :- Y/ a(X, Y); Z/ b(X, Z).
Y/ p(X, Y), Z/ q(X, Z) :- a(X).
Y/ p(X, Y); Z/ q(X, Z) :- a(X).
```

(1) relation determining a concept by composition:
```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

**conceptual merge by morphism**
*(take the prototype dog and make it angrier).* in logic usually operationalized as intersection.

(1) relation defining new concepts:

```
angrydog(X) :- dog(X), angry(X).
```

(2) relation defining taxonomical relations:

```
mammal(X) :- dog(X); cat(X).
```

(3) relation activating back the source concepts:

```
dog(X), angry(X) :- angrydog(X).
```

(4) non-deterministic relation activating alternative choices:

```
dog(X); cat(X) :- mammal(X).
```

**conceptual merge by aggregation** *(construct a whole out of components).*

(1) relation determining a concept by composition:

```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```

(2) specifies an operation of conceptual generalization:

```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```

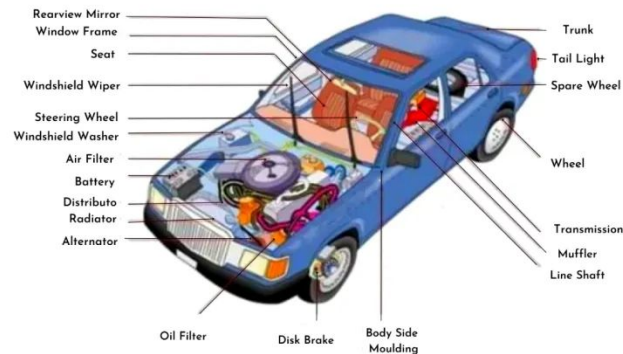(3) reifies mereonomical relations (similarly to a contructor in OOP):

```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```

(4) activates the possible realizations of a concept:

```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

**conceptual merge by morphism** *(take the prototype dog and make it angrier).* in logic usually operationalized as intersection.

(1) relation defining new concepts:

```
angrydog(X) :- dog(X), angry(X).
```

(2) relation defining taxonomical relations:

```
mammal(X) :- dog(X); cat(X).
```

(3) relation activating back the source concepts:

```
dog(X), angry(X) :- angrydog(X).
```

(4) non-deterministic relation activating alternative choices:

```
dog(X); cat(X) :- mammal(X).
```

**COMPREHENSION**

**conceptual merge by aggregation** *(construct a whole out of components).*

(1) relation determining a concept by composition:

```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```

(2) specifies an operation of conceptual generalization:

```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```

(3) reifies mereonomical relations (similarly to a contructor in OOP):

```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```

(4) activates the possible realizations of a concept:

```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

(1) relation defining new concepts:
```
angrydog(X) :- dog(X), angry(X).
```
(2) relation defining taxonomical relations:
```
mammal(X) :- dog(X); cat(X).
```
(3) relation activating back the source concepts:
```
dog(X), angry(X) :- angrydog(X).
```
(4) non-deterministic relation activating alternative choices:
```
dog(X); cat(X) :- mammal(X).
```

**GENERALIZATION**

(1) relation determining a concept by composition:
```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

(1) relation defining new concepts:
```
angrydog(X) :- dog(X), angry(X)
```

(2) relation defining taxonomical relations:
```
mammal(X) :- dog(X); cat(X).
```

(3) relation activating back the source concepts:
```
dog(X), angry(X) :- angrydog(X).
```

(4) non-deterministic relation activating alternative choices:
```
dog(X); cat(X) :- mammal(X).
```

**fusion**

**GENERALIZATION**

(1) relation determining a concept by composition:
```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```

(2) specifies an operation of conceptual generalization:
```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```

(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```

(4) activates the possible realizations of a concept:
```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

(1) relation defining new concepts:
```
angrydog(X) :- dog(X), angry(X).
```
(2) relation defining taxonomical relations:
```
mammal(X) :- dog(X); cat(X).
```
(3) relation activating back the source concepts:
```
dog(X), angry(X) :- angrydog(X).
```
(4) non-deterministic relation activating alternative choices:
```
dog(X); cat(X) :- mammal(X).
```

**GENERALIZATION**

(1) relation determining a concept by composition:
```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
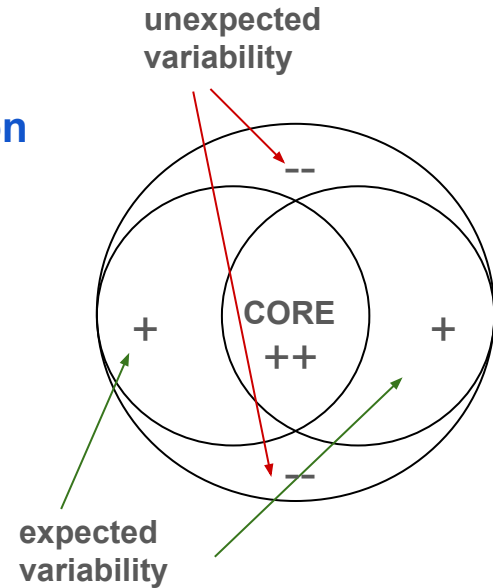(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```

**fusion**

**unexpected variability**

--

**CORE**
**++**

+

+

--

**expected variability**

(1) relation defining new concepts:
    angrydog(X) :- dog(X), angry(

(2) relation defining taxonomical relations:
    mammal(X) :- dog(X); cat(X).

(3) relation activating back the source concepts:
    dog(X), angry(X) :- angrydog(X).

(4) non-deterministic relation activating alternative choices:
    dog(X); cat(X) :- mammal(X).

**DESCRIPTION**



(1) relation determining a concept by composition:
    car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X,

(2) specifies an operation of conceptual generalization:
    student(X) :- Y/ humanities(Y), studies(X, Y); Z/ scienc            , Z).

(3) reifies mereonomical relations (similarly to a contructor in OOP):
    Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).

(4) activates the possible realizations of a concept:
    Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).

(1) relation defining new concepts:
    angrydog(X) :- dog(X), angry(

(2) relation defining taxonomical relations:
    mammal(X) :- dog(X); cat(X).

(3) relation activating back the source concepts:
    dog(X), angry(X) :- angrydog(X).

(4) non-deterministic relation activating alternative choices:
    dog(X); cat(X) :- mammal(X).

telling at large?

# DESCRIPTION

(1) relation determining a concept by composition:
    car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X,
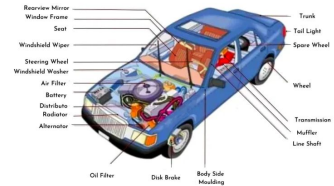
(2) specifies an operation of conceptual generalization:
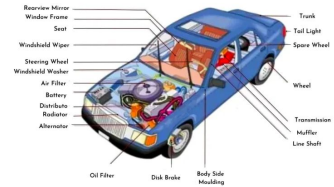    student(X) :- Y/ humanities(Y), studies(X, Y); Z/ scien                , Z).

(3) reifies mereonomical relations (similarly to a contructor in OOP):
    Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).

(4) activates the possible realizations of a concept:
    Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).

(1) relation defining new concepts:
```
angrydog(X) :- dog(X), angry(X).
```
(2) relation defining taxonomical relations:
```
mammal(X) :- dog(X); cat(X).
```
(3) relation activating back the source concept:
```
dog(X), angry(X) :- angrydog(X).
```
(4) non-deterministic relation activating alternative choices:
```
dog(X); cat(X) :- mammal(X).
```

**SPECIFICATION**





(1) relation determining a concept by composition:
```
car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
```
(2) specifies an operation of conceptual generalization:
```
student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
```
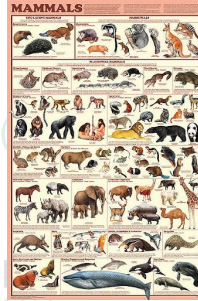(3) reifies mereonomical relations (similarly to a contructor in OOP):
```
Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).
```
(4) activates the possible realizations of a concept:
```
Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).
```
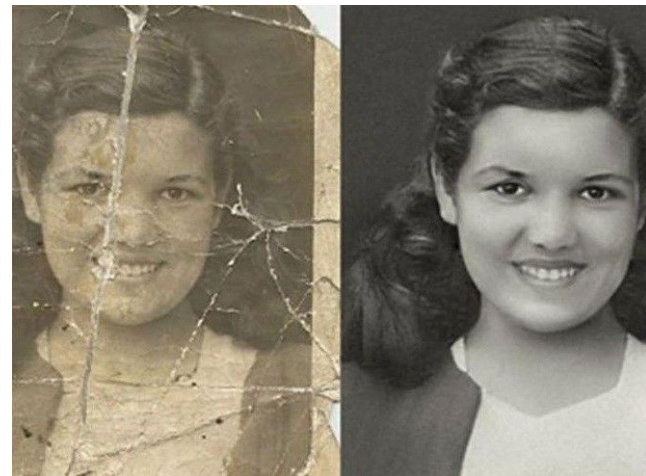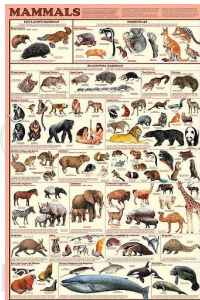
(1) relation defining new concepts:
    angrydog(X) :- dog(X), angry(X).
(2) relation defining taxonomical relations:
    mammal(X) :- dog(X); cat(X).
(3) relation activating back the source concept:
    dog(X), angry(X) :- angrydog(X).
(4) non-deterministic relation activating alternative choices:
    dog(X); cat(X) :- mammal(X).

"filling the gaps" mechanism
(including text completion, as in LLM)

(1) relation determining a concept by composition:
    car(X) :- Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z).
(2) specifies an operation of conceptual generalization:
    student(X) :- Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z).
(3) reifies mereonomical relations (similarly to a contructor in OOP):
    Y/ engine(Y), Z/ wheels(Z), has(X, Y), has(X, Z) :- car(X).

**SPECIFICATION**

(4) activates the possible realizations of a concept:
    Y/ humanities(Y), studies(X, Y); Z/ sciences(Z), studies(X, Z) :- student(X).

# A suite of mechanisms!

via *contrast*

**3. DESCRIPTION**

**un-packing**

vs

**packing**

**1. COMPREHENSION**

via *merge*

via *detachment*

**4. SPECIFICATION**

**zooming-in**

vs

**zooming-out**

**2. GENERALIZATION**

via *fusion*

# A suite of mechanisms!

decompression

via *contrast*

**3. DESCRIPTION**

**un-packing**

vs

**packing**

**1. COMPREHENSION**

via *merge*

via *detachment*

**4. SPECIFICATION**

**zooming-in**

vs

**zooming-out**

**2. GENERALIZATION**

via *fusion*

compression

# A suite of mechanisms!

decompression

via *contrast*                               via *detachment*

**3. DESCRIPTION**                          **4. SPECIFICATION**

**un-packing**                                  **zooming-in**

vs                                                      vs

**packing**                                     **zooming-out**

**1. COMPREHENSION**                     **2. GENERALIZATION**

via *merge*                                        via *fusion*

*re-encoding*          compression          *quantization*

# Conclusion

- The first mathematical model of a neuron was the McCulloch-Pitts (MCP) model in 1943, which paved the way to the Perceptron and then… to transformers (including LLMs).
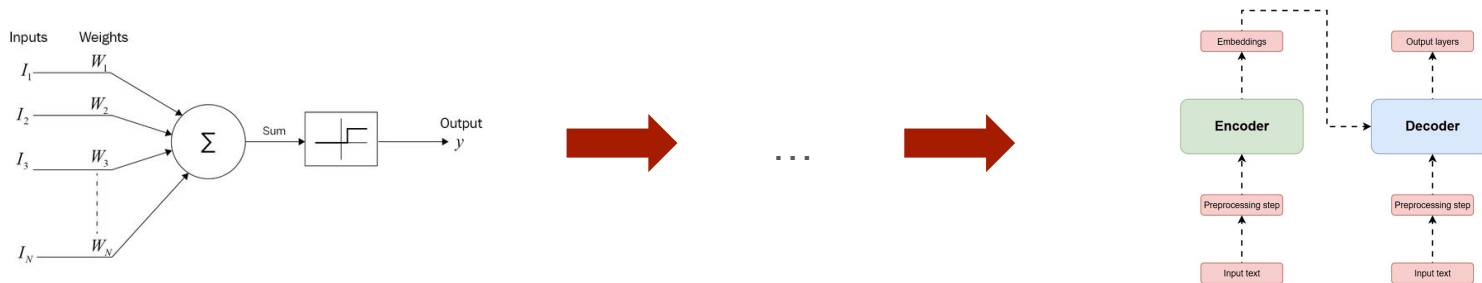
# Conclusion

- The first mathematical model of a neuron was the McCulloch-Pitts (MCP) model in 1943, which paved the way to the Perceptron and then… to transformers (including LLMs).



and now we're trying to understand what they do!

# Conclusion

-



```
1  % neutrality(+Matrix,+Exprs,-Exprs): the function 𝒩(X)
2  neutrality(AttM, X, Y) :-
3      mv_mult(AttM, X, Z),        % ℛ⁺(X)
4      maplist(bnot, Z, Y).
5
6  % innocuousity(+Matrix,+Exprs,-Exprs): the function ℐ(X)
7  innocuousity(AttM, X, Y) :-
8      transpose(AttM, AttM_t),   % transpose operation
9      mv_mult(AttM_t, X, Z).     % ℛ⁻(X)
10     maplist(bnot, Z, Y).
11
12 % defense(+Matrix,+Exprs,-Exprs): the function ℱ(X)
13 defense(AttM, X, Y) :-
14     neutrality(AttM, X, Z),
15     neutrality(AttM, Z, Y).
```

- In this paper, rather than trying to make sense of inferential circuits embedded in transformers, we looked at known inferential constructs by the lens of simple digital circuits…
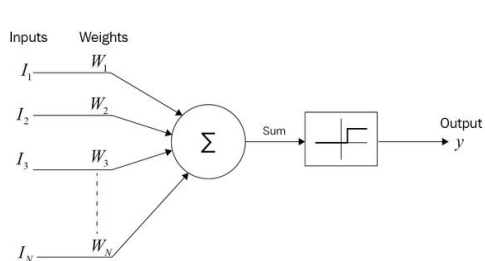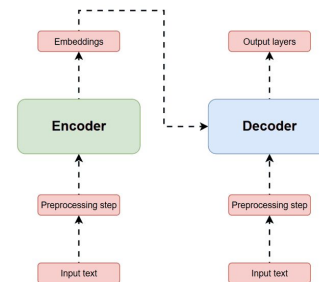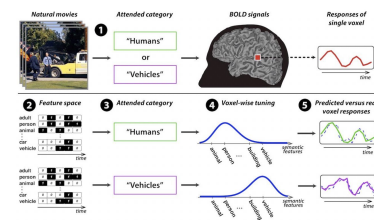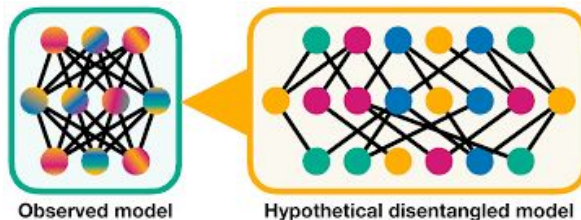
# Conclusion

- The first mathematical model of a neuron was the McCulloch-Pitts (MCP) model in 1943, which paved the way to the Perceptron and then… to transformers (including LLMs).
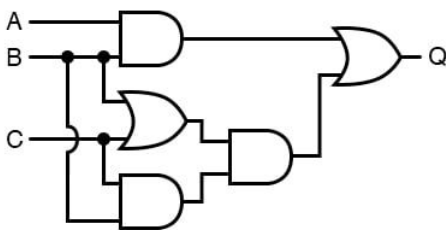


```
1  % neutrality(+Matrix,+Exprs,-Exprs): the function 𝒩(X)
2  neutrality(AttM, X, Y) :-
3      mv_mult(AttM, X, Z),       % ℛ⁺(X)
4      maplist(bnot, Z, Y).
5
6  % innocuousity(+Matrix,+Exprs,-Exprs): the function ℐ(X)
7  innocuousity(AttM, X, Y) :-
8      transpose(AttM, AttM_t),   % transpose operation
9      mv_mult(AttM_t, X, Z).     % ℛ⁻(X)
10     maplist(bnot, Z, Y).
11
12 % defense(+Matrix,+Exprs,-Exprs): the function ℱ(X)
13 defense(AttM, X, Y) :-
14     neutrality(AttM, X, Z),
15     neutrality(AttM, Z, Y).
```

- In this paper, rather than trying to make sense of inferential circuits embedded in transformers, we looked at known inferential constructs by the lens of simple digital circuits… and we found that:

    **there exists a barebone common structure for inferential mechanisms!**

# Further results!

- In the paper we have additional results:

    - an application of a probabilistic interpretation of logic programs predicting dependencies across the 4 inferential mechanisms,

    - an elaboration on how learning would work differently for generalization and comprehension,

    - an observation on the fact that LLMs mimic a description mechanism through a specification one, entailing that they may take contrast wrong.

    - and yet others considerations!!

# Exploring structures of inferential mechanisms through simplistic digital circuits

25 October 2025, AIC workshop @ ECAI 2025
*10th International Workshop on Artificial Intelligence and Cognition*

Giovanni Sileno
g.sileno@uva.nl

*University of Amsterdam*

Jean-Louis Dessalles
jean-louis.dessalles@telecom-paris.fr

*Télécom Paris*