

Université de Montréal - Université Côte d'Azur
Colloque Créativités Artificielles
Nice, France

Critique and Formalism

The Role of Formal Sciences for an Internalist Critique of LLMs

Juan Luis Gastaldi

`www.giannigastaldi.com`

ETH zürich

April 28, 2026

Introduction: Critique and Formalism

Empiricism and Formalism in Machine Learning

Formal Explainability

Theoretical Interpretability

Conclusion: AI and Structuralism

Introduction: Critique and Formalism

Empiricism and Formalism in Machine Learning

Formal Explainability

Theoretical Interpretability

Conclusion: AI and Structuralism

Where Art Thou, Critique?

- ◇ Kirschenbaum (2023):
Bender et al.'s (2021) paper “offers a **disarmingly linear account of how language, communication, intention, and meaning work**, one that would seem to sidestep decades of scholarship around these same issues in literary theory [...] the passage would be **red meat for a graduate critical-theory seminar.**”
- ◇ Underwood (2023):
“The beautiful **irony** of this situation [...] is that a generation of **humanists trained on Foucault** have now rallied around “On the Dangers of Stochastic Parrots” to **oppose a theory of language that their own disciplines invented**, just at the moment when computer scientists are reluctantly beginning to accept it.”

Introduction: Critique and Formalism

Empiricism and Formalism in Machine Learning

Formal Explainability

Theoretical Interpretability

Conclusion: AI and Structuralism

An Empiricist Turn?

Interpretability as a Natural Science

The Structure of Scientific Revolutions by Thomas Kuhn [42] is a classic text on the history and sociology of science. In it, Kuhn distinguishes between “normal science” in which a scientific community has a paradigm, and “extraordinary science” in which a community lacks a paradigm, either because it never had one or because it was weakened by crisis. It’s worth noting that “extraordinary science” is not a desirable state: it’s a period where researchers struggle to be productive.

Kuhn’s description of pre-paradigmatic fields feel eerily reminiscent of interpretability today. ⁹ There isn’t consensus on what the objects of study are, what methods we should use to answer them, or how to evaluate research results. To quote a recent interview with Ian Goodfellow: “For interpretability, I don’t think we even have the right definitions.” [43]

One particularly challenging shared sense of how to deal with this, especially those with which can evaluate HCI background may

But interpretability of neural networks are biology. Such work will be held to the standard

The Empiricization of Computer Science



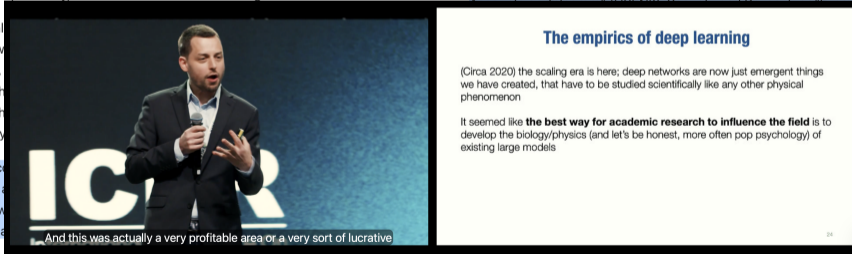
MANOEL HORTA RIBEIRO

DEC 17, 2025



Share

Yet the nature of computer science has changed tremendously: it is now an empirical science, driven by observation and experiment as much as by theory and construction. If you don’t believe me, spend 5 minutes going over the best paper awards for three conferences in different subfields in CS. If you did so in 2025, you might have found papers like “Characterizing and Detecting Propaganda-Spreading in Social Media” (ICSEW), “Evaluating the Impact of AI on Mental Health” (CHI), and “Human-Centered AI for Education” (CHI). These papers describe online user behavior, data, people using their own examination of social networks. They are published independently



(Olah et al., 2020)

Zico Kolter, *Building Safe and Robust AI Systems*, Keynote at ICLR 2025.

(Horta Ribeiro, 2025)

$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$

$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$$

$$0: \lambda f. \lambda x. x$$

$$1: \lambda f. \lambda x. f x$$

$$2: \lambda f. \lambda x. f (f x)$$

$$3: \lambda f. \lambda x. f (f (f x))$$

$$4: \lambda f. \lambda x. f (f (f (f x)))$$

$$5: \lambda f. \lambda x. f (f (f (f (f x))))$$

$$\dots$$

$$n: \lambda f. \lambda x. \underbrace{f(\dots(f x)\dots)}_{n \text{ times}}$$

$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$$

$$0: \lambda f. \lambda x. x$$

$$\lambda m. \lambda n. \lambda f. \lambda x. m f (n f x) (\lambda f. \lambda x. f (f x)) (\lambda f. \lambda x. f (f (f x)))$$

$$1: \lambda f. \lambda x. f x$$

$$2: \lambda f. \lambda x. f (f x)$$

$$3: \lambda f. \lambda x. f (f (f x))$$

$$4: \lambda f. \lambda x. f (f (f (f x)))$$

$$5: \lambda f. \lambda x. f (f (f (f (f x))))$$

$$\dots$$

$$n: \lambda f. \lambda x. \underbrace{f(\dots(f x)\dots)}_{n \text{ times}}$$

$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$

0: $\lambda f. \lambda x. x$

1: $\lambda f. \lambda x. f x$

2: $\lambda f. \lambda x. f (f x)$

3: $\lambda f. \lambda x. f (f (f x))$

4: $\lambda f. \lambda x. f (f (f (f x)))$

5: $\lambda f. \lambda x. f (f (f (f (f x))))$

...

$n: \lambda f. \lambda x. \underbrace{f(\dots(f x)\dots)}_{n \text{ times}}$

$\lambda m. \lambda n. \lambda f. \lambda x. m f (n f x) (\lambda f. \lambda x. f (f x)) (\lambda f. \lambda x. f (f (f x)))$

⋮

⋮

⋮

⋮

⋮

⋮

⋮

$\lambda f. \lambda x. f (f (f (f (f x))))$

$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$

$P' := \lambda r. \lambda s. \lambda f. \lambda x. f(f(f(f(fx))))$

0: $\lambda f. \lambda x. x$

1: $\lambda f. \lambda x. f x$

2: $\lambda f. \lambda x. f(fx)$

3: $\lambda f. \lambda x. f(f(fx))$

4: $\lambda f. \lambda x. f(f(f(fx)))$

5: $\lambda f. \lambda x. f(f(f(f(fx))))$

...

$n: \lambda f. \lambda x. \underbrace{f(\dots(fx)\dots)}_{n \text{ times}}$

$\lambda r. \lambda s. \lambda f. \lambda x. f(f(f(f(fx))))(\lambda f. \lambda x. f(fx))(\lambda f. \lambda x. f(f(fx)))$

↘

↘

↘

↘

↘

↘

↘

$\lambda f. \lambda x. f(f(f(f(fx))))$

$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$$

$$0: \lambda f. \lambda x. x$$

$$1: \lambda f. \lambda x. f x$$

$$2: \lambda f. \lambda x. f (f x)$$

$$3: \lambda f. \lambda x. f (f (f x))$$

$$4: \lambda f. \lambda x. f (f (f (f x)))$$

$$5: \lambda f. \lambda x. f (f (f (f (f x))))$$

...

$$n: \lambda f. \lambda x. \underbrace{f(\dots(f x)\dots)}_{n \text{ times}}$$

$$\lambda m. \lambda n. \lambda f. \lambda x. m f (n f x) (\lambda f. \lambda x. f (f x)) (\lambda f. \lambda x. f (f (f x)))$$

⋮

⋮

⋮

⋮

⋮

⋮

⋮

$$\lambda f. \lambda x. f (f (f (f (f x))))$$

$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$$

$$0: \lambda f. \lambda x. x$$

$$1: \lambda f. \lambda x. f x$$

$$2: \lambda f. \lambda x. f (f x)$$

$$3: \lambda f. \lambda x. f (f (f x))$$

$$4: \lambda f. \lambda x. f (f (f (f x)))$$

$$5: \lambda f. \lambda x. f (f (f (f (f x))))$$

...

$$n: \lambda f. \lambda x. \underbrace{f(\dots(f x)\dots)}_{n \text{ times}}$$

$$\lambda m. \lambda n. \lambda f. \lambda x. m f (n f x) (\lambda f. \lambda x. f (f x)) (\lambda f. \lambda x. f (f (f x)))$$

$$\lambda m. \lambda n. \lambda f. \lambda x. m f (n f x) (\lambda g. \lambda y. g (g y)) (\lambda h. \lambda z. h (h (h z)))$$

$$\lambda n. \lambda f. \lambda x. (\lambda g. \lambda y. g (g y)) f (n f x) (\lambda h. \lambda z. h (h (h z)))$$

$$\lambda n. \lambda f. \lambda x. (\lambda g. \lambda y. g (g y)) f (n f x) (\lambda h. \lambda z. h (h (h z)))$$

$$\lambda f. \lambda x. (\lambda g. \lambda y. g (g y)) f ((\lambda h. \lambda z. h (h (h z)))) f x$$

$$\lambda f. \lambda x. (\lambda y. f (f y)) ((\lambda h. \lambda z. h (h (h z)))) f x$$

$$\lambda f. \lambda x. (\lambda y. f (f y)) ((\lambda z. f (f (f z)))) x$$

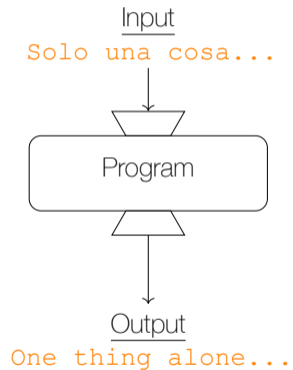
$$\lambda f. \lambda x. (\lambda y. f (f y)) (f (f (f x)))$$

$$\lambda f. \lambda x. f (f (f (f (f x))))$$

$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$

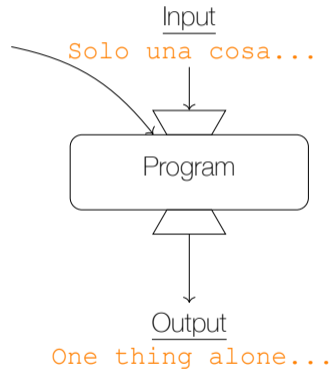
$P'' := \lambda R \acute{o} f \ddot{A} \ddot{O} \hat{e} \ddot{N} 5 \ddot{E} | \ddot{A} x \ddot{n} = \infty \ddot{u} \ddot{y} m W f 286 \ddot{e} y ' S \ddot{O} \acute{u} > v \& i \hat{A} - 2 \acute{o} \acute{E} 7 \acute{o} \zeta \infty \{ \ddot{a} > 2 f \ddot{B} \acute{o} \mu G \# \ddot{A} 9 \zeta U$
 $\infty \text{ob} t Y B \hat{o} \ddot{Y} \ddot{U} \ddot{e} \% 0 3 ; 5 \acute{a} [l - \acute{e} u \hat{o} \ddot{U} \acute{e} \acute{7} - \ddot{U} . \lambda : \hat{^} 4 m \acute{O} \acute{O} \ddot{Y} ' \acute{e} - + \acute{I} s \acute{O} , \$ + g \ddot{i} , B^{TM} \div o - \# i \ddot{Y} \hat{e} \hat{U} v$
 $- g \acute{O} \ddot{y} / \acute{e} i i j O \ddot{f} \acute{C} e f i \bullet J 1 \ll \acute{E} \acute{o} , \acute{I} h \hat{e} t \ddot{f} \acute{a} e Y \$ \hat{^} 6 F i W \gg R \acute{U} K g e \ddot{r} . \lambda \ddot{f} d \ddot{r} \dots D 2 \div \acute{c} \acute{o} \acute{x} \acute{e} \acute{E} y . \acute{O} \ddot{r} c b$
 $B \acute{e} \acute{L} N \acute{E} 1 \hat{E} \ddot{f} / \hat{U} 9 \ddot{N} \mu - / J Y \zeta \acute{o} \acute{E} 9 \ddot{y} \hat{A} \acute{E} \acute{E} . \lambda \acute{A} \acute{I} \hat{A} \hat{^} \acute{o} \zeta , \gg f q \infty \pm \hat{i} \sim B 5 \hat{I} > O \sim g^{TM} \acute{6} \Omega e \acute{a} \acute{e} C / \acute{a} \dots \acute{O}$
 $\cdot f \acute{O} \acute{A}] \ddot{N} \acute{a} y \hat{E} N \acute{e} \hat{E} \ddot{r} . \lambda \acute{A} \acute{e} \acute{a} \acute{e} f U \acute{o} \acute{f} E \acute{U} \acute{I} m \# , , 4 \backslash r \sqrt{-} \div \hat{I} p \acute{o} \gg y \ast v t \acute{A} J \acute{A} F 1 \hat{u} \acute{A} \acute{o} z \ll \acute{n} M \ddot{r} D j \acute{C} E$
 $B \acute{E} \acute{e} \acute{I} T _ \hat{E} a \% 0 \acute{A} \zeta \Omega @ \backslash \acute{O} \hat{^} \sim] \hat{I} \ddot{h} \ddot{f} : \hat{^} 4 m \acute{O} \acute{O} \ddot{Y} ' \acute{e} - + \acute{I} s \acute{O} , \$ + g \ddot{i} , B^{TM} \div o - \# i \ddot{Y} \hat{e} \hat{U} v - g \acute{O} \ddot{y}$
 $/ \acute{e} i i j O \ddot{f} \acute{C} e f i \bullet J 1 \ll \acute{E} \acute{o} , \acute{I} h \hat{e} t \ddot{f} \acute{a} e Y \$ \hat{^} 6 F i W \gg R \acute{U} K g e \ddot{r} \acute{A} \acute{I} \hat{A} \hat{^} \acute{o} \zeta , \gg f q \infty \pm \hat{i} \sim B 5 \hat{I} > O \sim g^{TM} \acute{6}$
 $\Omega e \acute{a} \acute{e} C / \acute{a} \dots \acute{O} \cdot f \acute{O} \acute{A}] \ddot{N} \acute{a} y \hat{E} N \acute{e} \hat{E} \ddot{r} (\ddot{f} d \ddot{r} \dots D 2 \div \acute{c} \acute{o} \acute{x} \acute{e} \acute{E} y . \acute{O} \ddot{r} c b B \acute{e} \acute{L} N \acute{E} 1 \hat{E} \ddot{f} / \hat{U} 9 \ddot{N} \mu - /$
 $J Y \zeta \acute{o} \acute{E} 9 \ddot{y} \hat{A} \acute{E} \acute{E} \acute{A} \acute{I} \hat{A} \hat{^} \acute{o} \zeta , \gg f q \infty \pm \hat{i} \sim B 5 \hat{I} > O \sim g^{TM} \acute{6} \Omega e \acute{a} \acute{e} C / \acute{a} \dots \acute{O} \cdot f \acute{O} \acute{A}] \ddot{N} \acute{a} y \hat{E} N \acute{e} \hat{E} \ddot{r} \acute{A}$
 $\acute{e} \acute{f} U \acute{o} \acute{f} E \acute{U} \acute{I} m \# , , 4 \backslash r \sqrt{-} \div \hat{I} p \acute{o} \gg y \ast v t \acute{A} J \acute{A} F 1 \hat{u} \acute{A} \acute{o} z \ll \acute{n} M \ddot{r} D j \acute{C} E B \acute{E} \acute{e} \acute{I} T _ \hat{E} a \% 0 \acute{A} \zeta \Omega @ \backslash$
 $\acute{O} \hat{^} \sim] \hat{I} \ddot{h} \ddot{f}) (\acute{E} \hat{I} \hat{U} \acute{e} i 4 W \mu \acute{I} \} w , , \$ \Omega \acute{K} 5 \acute{e} \acute{A} \acute{Q} \% 3 [m \acute{r} \sim B \acute{A} f i \acute{f} \acute{O} ; \acute{o} J \zeta C \acute{E} \hat{i} \acute{o} \ddot{Y} \acute{O} c B , \acute{n} \$ \acute{A} \acute{a} \} \acute{O} \acute{A} \acute{O} 3 ;$
 $\sim ? \acute{o} \acute{o} \acute{C} E @ f \acute{l} 8 \sim R C \acute{A} \acute{e} \acute{o} \sim \ast \& < \acute{Y} - \acute{o} 1 2 \acute{A} \% 0 \acute{a} \acute{O} \acute{U} \# \acute{i} \acute{r} , \acute{u} \acute{r} \ll \acute{o} \acute{r} , , \infty \acute{I} \acute{a} \acute{a} \acute{e} \acute{O} \acute{A} d | \sim \acute{N} \acute{r} \acute{E} y \acute{O} ; \hat{^} W$
 $\ddot{r} \acute{w} \acute{o} [] \backslash \gg \acute{O} \acute{E} \acute{u} w \acute{r} 6 < \acute{u} \acute{r} = \acute{a} \acute{O} \acute{r} \acute{I} \acute{D} z ? 2 \pm | \acute{e} \acute{r} 3 \hat{A} / r x \mu \infty \mu \$ \acute{A} \acute{e} \hat{A} \ast l \acute{f} \sim \hat{i} \acute{u} \acute{r} + \acute{I} V \acute{i} y \acute{a} G \acute{a} \acute{e} \acute{B} \acute{a} g \acute{o} / , u \ddot{N}$

Implicit Structure



Implicit Structure

Algorithm
Translation

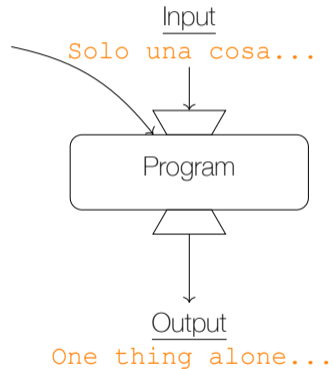


Implicit Structure

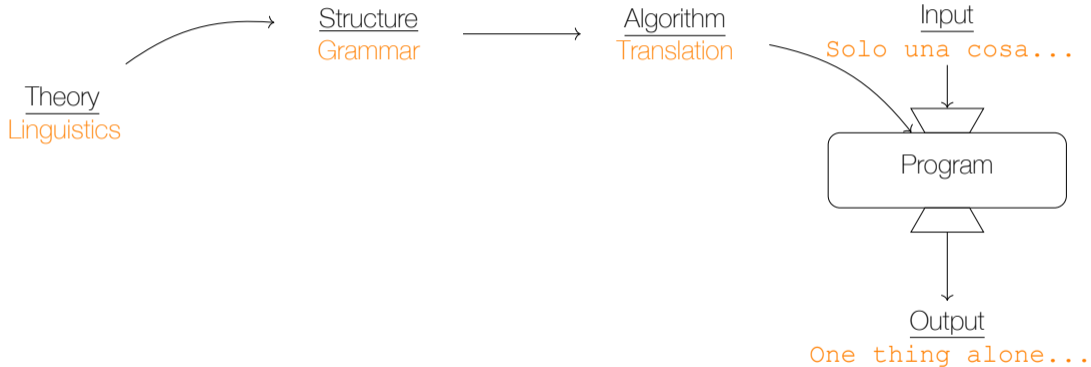
Structure
Grammar



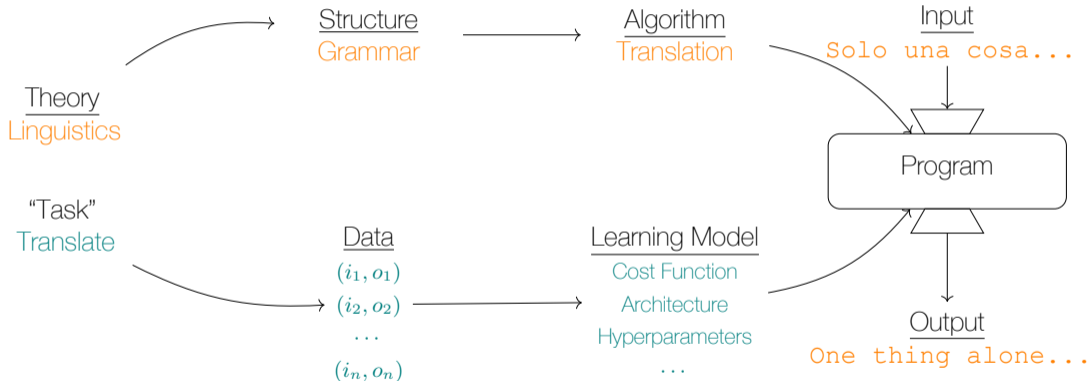
Algorithm
Translation



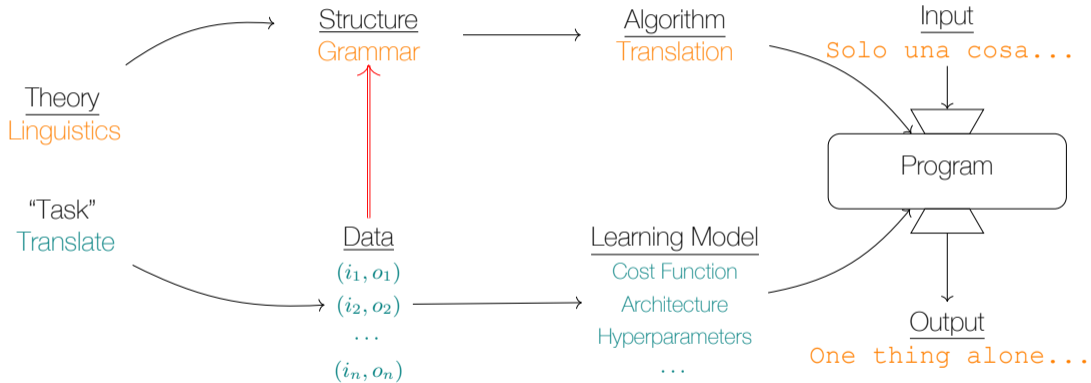
Implicit Structure



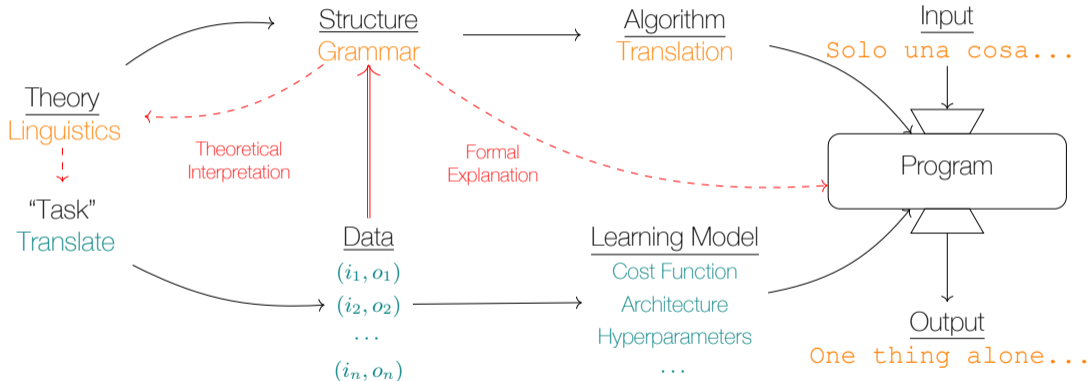
Implicit Structure



Implicit Structure



Making It Explicit



Introduction: Critique and Formalism

Empiricism and Formalism in Machine Learning

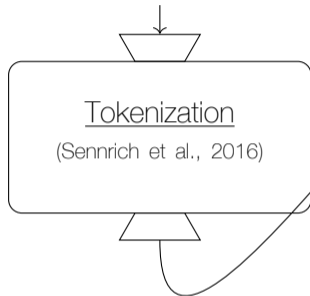
Formal Explainability

Theoretical Interpretability

Conclusion: AI and Structuralism

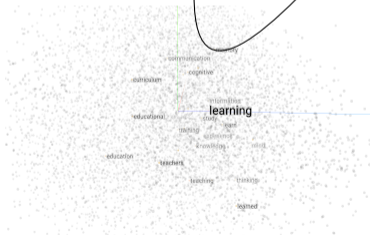
Embeddings in LLMs

Epistemology of Machine Learning
Distributional Language Models

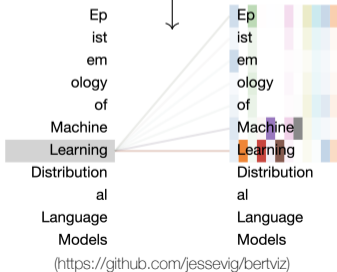
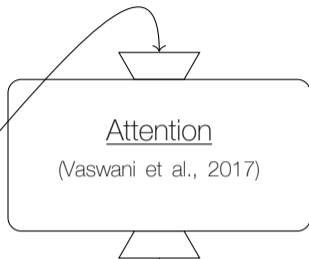


Epistemology of Machine Learning
Distributional Language Models

(<https://tiktokenizer.vercel.app>)



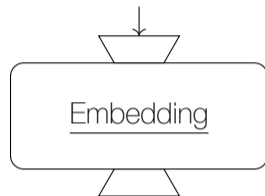
(<https://projector.tensorflow.org>)



(<https://github.com/jessevig/bertviz>)

Neural Embeddings

Epistemology of Machine Learning
Distributional Language Models



Neural Embeddings

Structure

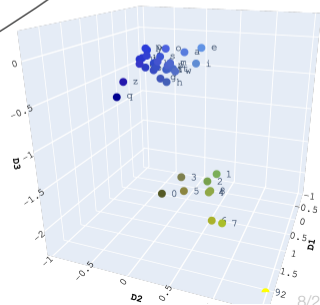
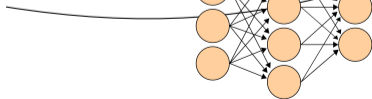


{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}



Embedding

Data



$$\ell = \sum_{w \in V_w} \sum_{c \in V_c} \#(w, c) (\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)])$$

$$\frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{when} \quad \vec{w} \cdot \vec{c} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

- ◊ Word2vec performs an **implicit, low-dimensional factorization** of a **pointwise mutual information (pmi), word-context matrix**.
- ◊ The **Singular Value Decomposition (SVD)** provides an **exact solution** to this problem.

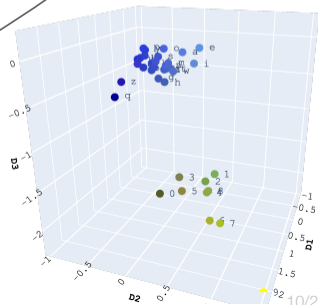
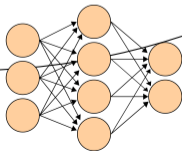
Embedding structure

Structure

?

$\{-, /, 0, 1, 2, \dots, 8, 9, =,$
 $a, b, c, \dots, w, x, y, z, \acute{e}\}$

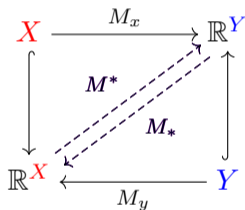
Embedding



From matrices to distributional operators

$$M: X \times Y \rightarrow \mathbb{R}$$

$$(x, y) \mapsto \text{pmi}(x, y)$$



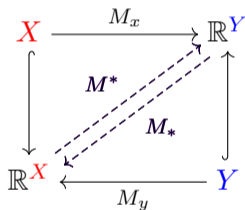
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$

From matrices to distributional operators

$$M: X \times Y \rightarrow \mathbb{R}$$

$$(x, y) \mapsto \text{pmi}(x, y)$$



$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$

$$\{u_1, \dots, u_m\} \subset \mathbb{R}^X$$

$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

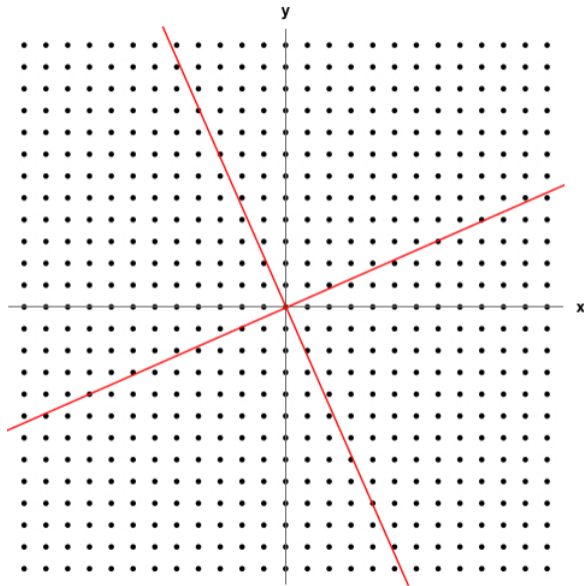
$$\{\lambda_1, \dots, \lambda_{\min(m,n)}, 0, \dots, 0\}$$

$$M_* M^* u_i = \lambda_i u_i$$

$$M^* M_* v_i = \lambda_i v_i$$

The u_i and v_i are (linear)
fixed points!

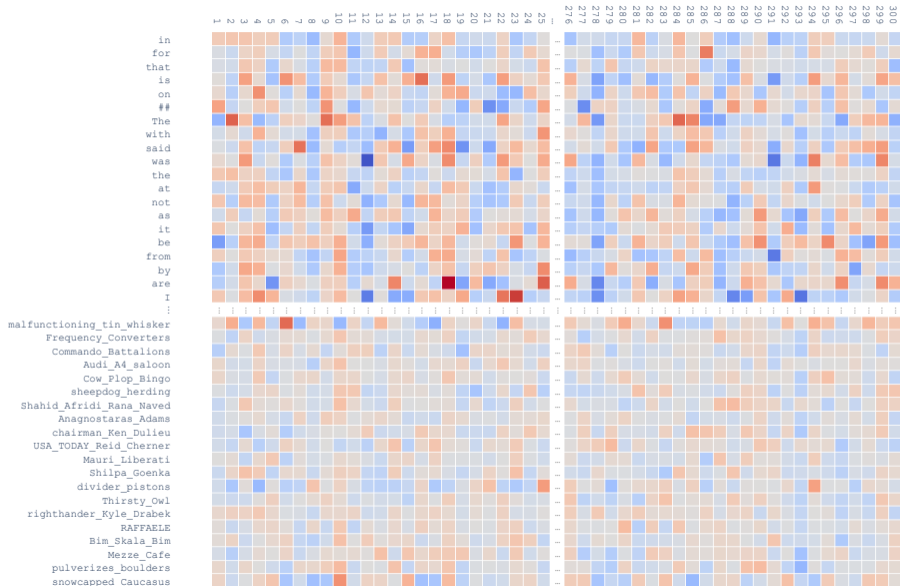
Eigenvectors



Embedding dimensions as fixed points



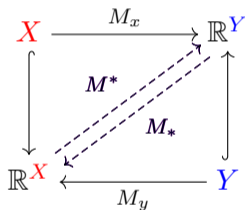
Embedding dimensions as fixed points



Embeddings as Functors Over Categories

$$M: X \times Y \rightarrow \mathbb{R}$$

$$(x, y) \mapsto \text{pmi}(x, y)$$



$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

$$M^* M_* : \mathbb{R}^Y \rightarrow \mathbb{R}^Y$$

$$\{u_1, \dots, u_m\} \subset \mathbb{R}^X$$

$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

$$\{\lambda_1, \dots, \lambda_{\min(m,n)}, 0, \dots, 0\}$$

$$M_* M^* u_i = \lambda_i u_i$$

$$M^* M_* v_i = \lambda_i v_i$$

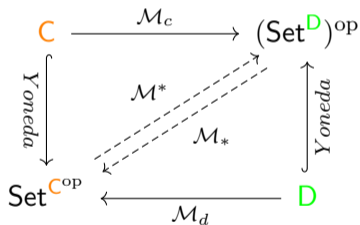
The u_i and v_i are (linear)
fixed points!

Embeddings as Functors Over Categories

$$\mathbf{C}^{\text{op}} \times \mathbf{D} \rightarrow \mathbf{Set}$$

$$\Downarrow$$

$$\mathcal{M}^* : \mathbf{Set}^{\mathbf{C}^{\text{op}}} \rightleftarrows (\mathbf{Set}^{\mathbf{D}})^{\text{op}} : \mathcal{M}_*$$



$$\mathcal{M}_* \mathcal{M}^* : \mathbf{Set}^{\mathbf{C}^{\text{op}}} \rightarrow \mathbf{Set}^{\mathbf{C}^{\text{op}}}$$

$$\mathcal{M}^* \mathcal{M}_* : (\mathbf{Set}^{\mathbf{D}})^{\text{op}} \rightarrow (\mathbf{Set}^{\mathbf{D}})^{\text{op}}$$

$$\text{Fix}(\mathcal{M}_* \mathcal{M}^*) := \{f \in \mathbf{Set}^{\mathbf{C}^{\text{op}}} \mid \mathcal{M}_* \mathcal{M}^*(f) \cong f\}$$

$$\text{Fix}(\mathcal{M}^* \mathcal{M}_*) := \{g \in (\mathbf{Set}^{\mathbf{D}})^{\text{op}} \mid \mathcal{M}^* \mathcal{M}_*(g) \cong g\}$$

Nucleus of $\mathcal{M} = \{(f_i, g_i)\}$, such that:

$$\mathcal{M}^* f_i \cong g_i \text{ and } \mathcal{M}_* g_i \cong f_i$$

The nucleus is a **category complete** and **cocomplete**

Categories \mathbf{C} and \mathbf{D} can be enriched!

$$\text{E.g.: } \mathcal{M}^* : \bar{\mathbf{R}}^{\mathbf{C}^{\text{op}}} \rightleftarrows (\bar{\mathbf{R}}^{\mathbf{D}})^{\text{op}} : \mathcal{M}_*$$

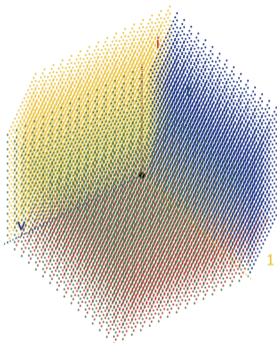
The Nucleus of the Profunctor

$$\begin{array}{ccc} e_i & s_i & \text{measurement} \\ \text{(terms)} & \text{(contexts)} & \\ \downarrow & \downarrow & \swarrow \\ C^{\text{op}} & \times D & \rightarrow \bar{\mathbb{R}} \end{array}$$

The Nucleus of the Profunctor

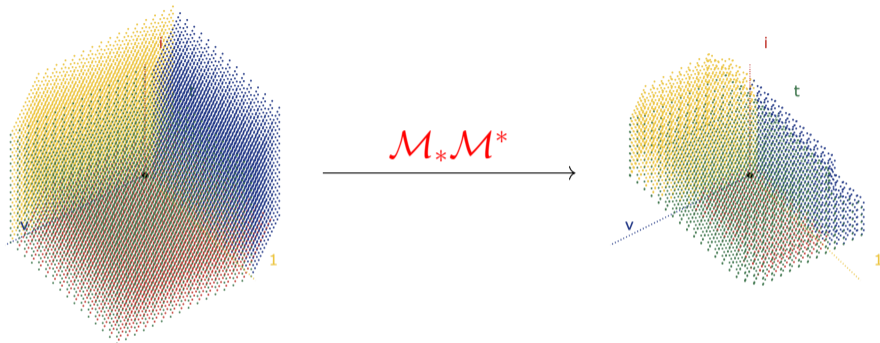
$$\begin{array}{ccc} e_i & s_i & \text{measurement} \\ \text{(terms)} & \text{(contexts)} & \\ \downarrow & \downarrow & \swarrow \\ \mathbf{C}^{\text{op}} \times \mathbf{D} & \rightarrow & \bar{\mathbb{R}} \\ \Downarrow & & \\ \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbf{C}^{\text{op}}} & \rightleftharpoons & (\bar{\mathbb{R}}^{\mathbf{D}})^{\text{op}} : \mathcal{M}_* \end{array}$$

The Nucleus of the Profunctor



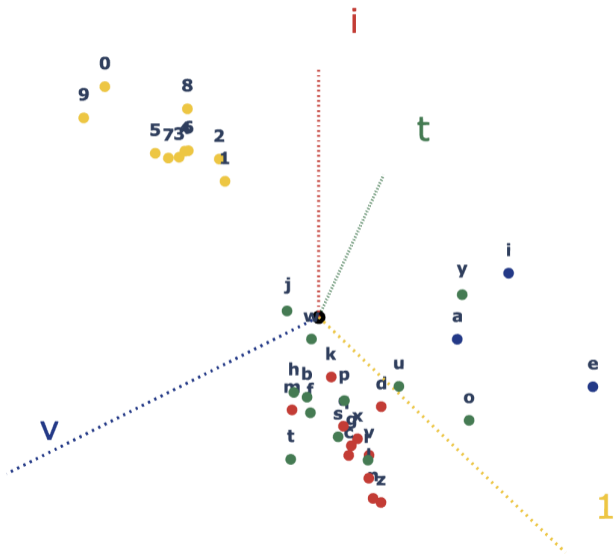
$$\begin{array}{ccc} \mathbb{C}^{\text{op}} \times \mathbb{D} & \rightarrow & \bar{\mathbb{R}} \\ & \Downarrow & \\ \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbb{C}^{\text{op}}} & \rightleftharpoons & (\bar{\mathbb{R}}^{\mathbb{D}})^{\text{op}} : \mathcal{M}_* \end{array}$$

The Nucleus of the Profunctor

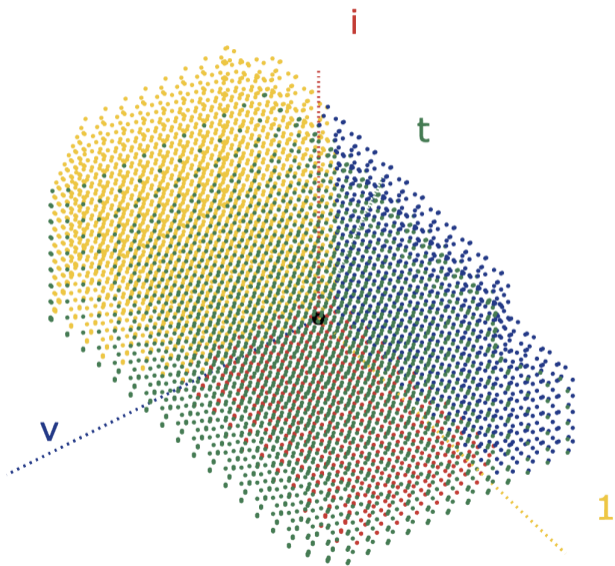


$$\begin{array}{ccc}
 \mathbb{C}^{\text{op}} \times \mathbb{D} & \rightarrow & \bar{\mathbb{R}} \\
 \Downarrow & & \\
 \mathcal{M}^* : \bar{\mathbb{R}}^{\mathbb{C}^{\text{op}}} & \rightleftharpoons & (\bar{\mathbb{R}}^{\mathbb{D}})^{\text{op}} : \mathcal{M}_*
 \end{array}$$

Geometry of the Nucleus



Geometry of the Nucleus



Introduction: Critique and Formalism

Empiricism and Formalism in Machine Learning

Formal Explainability

Theoretical Interpretability

Conclusion: AI and Structuralism

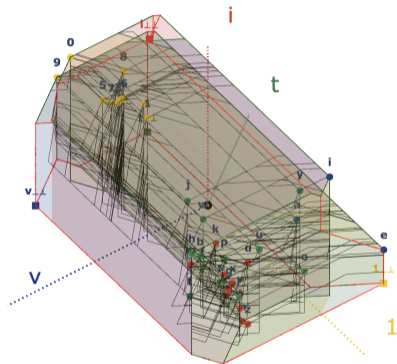
From the Distributional to the Structuralist Hypothesis

Theory
"Task"

?



Structure



From the Distributional to the Structuralist Hypothesis

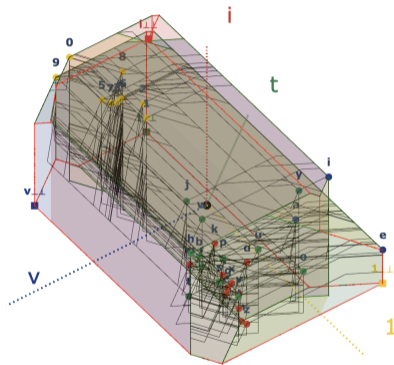
$$C^{\text{op}} \times D \rightarrow \bar{R}$$

Structure

Distributional Hypothesis

The content of linguistic units is determined by their *distribution* in a corpus.

Theory
"Task"



From the Distributional to the Structuralist Hypothesis

$$C^{op} \times D \rightarrow \bar{R}$$

Structure

Distributional Hypothesis

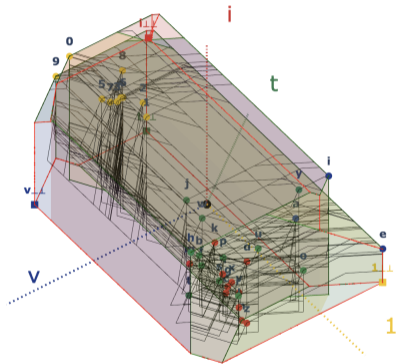
The content of linguistic units is determined by their *distribution* in a corpus.

Theory
"Task"



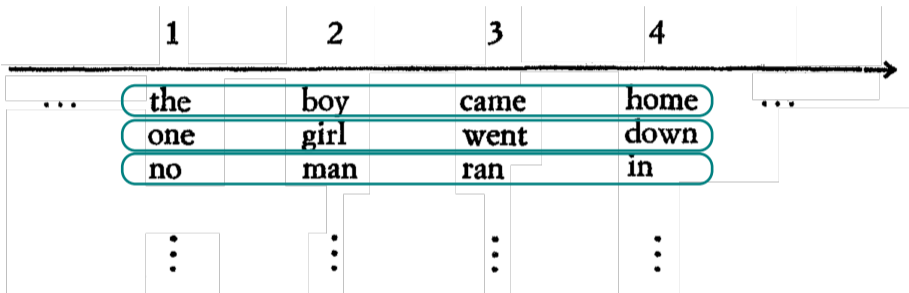
Structuralist Hypothesis

Linguistic content is the effect of a virtual *structure* underlying linguistic practices within a community

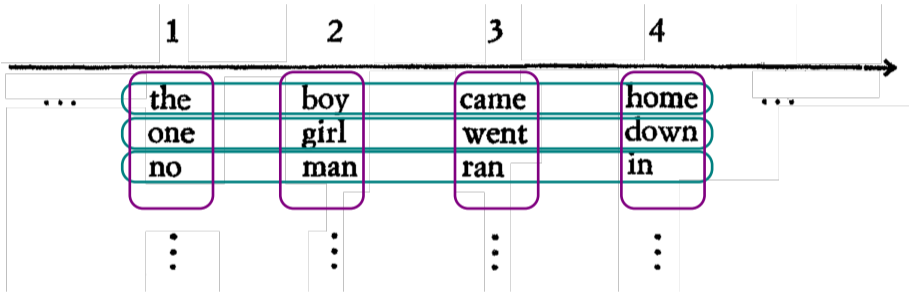


$$\bar{R}^{C^{op}} \Leftrightarrow (\bar{R}^D)^{op}$$

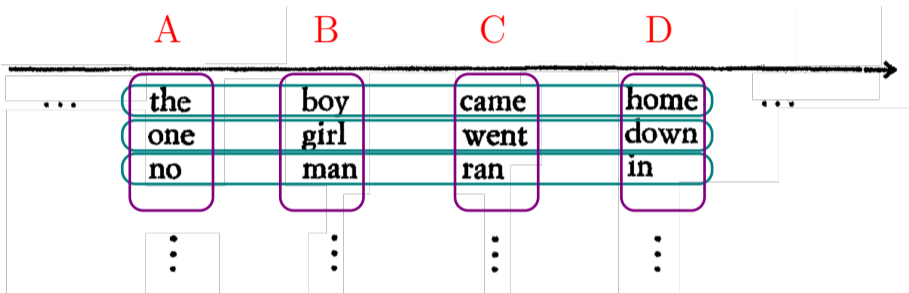
Syntagmas and Paradigmes



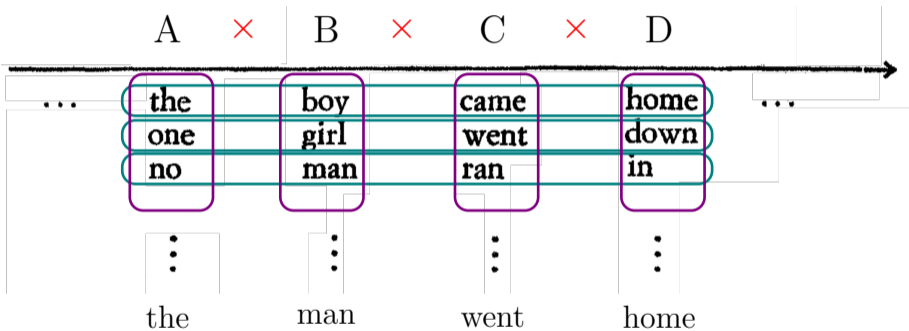
Syntagmas and Paradigmes



Syntagmas and Paradigmes



Syntagmas and Paradigmes



(Hjelmslev, 1971a)

Introduction: Critique and Formalism

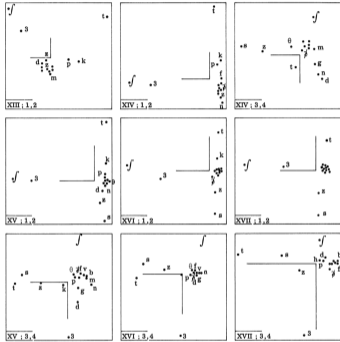
Empiricism and Formalism in Machine Learning

Formal Explainability

Theoretical Interpretability

Conclusion: AI and Structuralism

Structuralist Formalism



| v.p. | 1* | 2* | 3* | 4* | 5* |
|------|-----|-----|-----|-----|-----|
| XII | ,99 | ,98 | ,97 | ,92 | ,86 |
| XIII | ,94 | ,92 | ,86 | ,83 | ,76 |
| XIV | ,87 | ,62 | ,49 | ,47 | ,45 |
| XV | ,71 | ,41 | ,37 | ,27 | ,19 |
| XVI | ,55 | ,29 | ,25 | ,19 | ,04 |
| XVII | ,44 | ,23 | ,17 | ,09 | ,02 |

(benzécrid1976codage)

- 1) une voyelle neutre (amorphe), caractérisée par l'absence de chacune des propriétés $\beta, \phi, \chi, \lambda$: [a];
- 2) quatre types élémentaires de voyelles, chacun caractérisé par une seule propriété:

$$[e] = \phi, [a] = \beta, [v] = \chi, [z] = \lambda$$
- 3) six voyelles distinctes, chacune caractérisée par deux propriétés:

$$[o] = \phi\beta, [i] = \chi\phi, [u] = \chi\beta, [e] = \lambda\phi, [o] = \lambda\beta, [\partial] = \chi\lambda;$$
- 4) quatre voyelles combinées, chacune caractérisée par trois propriétés:

$$[e] = \chi\lambda\phi, [o] = \chi\lambda\beta, [u] = \phi\beta\chi, [\partial] = \phi\beta\lambda;$$
- 5) une voyelle polymorphe, caractérisée par les quatre propriétés considérées: la voyelle russe [sɨ] = $\phi\beta\chi\lambda$.

On obtient alors le diagramme de la figure 2.

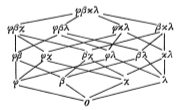
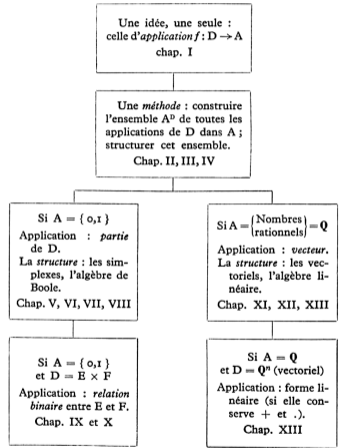


FIG. 2.

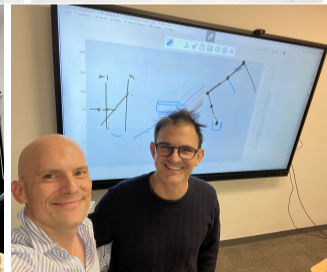
(Marcus, 1967)

Organigramme



(Barbut1967_T1)

Collaborations



J. Terilla (CUNY), T.-D. Bradley (SandboxAQ), L. Pellissier (Paris-Est Créteil), Th. Seiller (CNRS), S. Jarvis (CUNY)

Reference Papers

- ◇ Gastaldi, J. L. (2021). Why Can Computers Understand Natural Language? *Philosophy & Technology*, 34(1), 149–214. <https://doi.org/10.1007/s13347-020-00393-9>
- ◇ Gastaldi, J. L., & Pellissier, L. (2021). The Calculus of Language: Explicit Representation of Emergent Linguistic Structure through Type-Theoretical Paradigms. *Interdisciplinary Science Reviews*, 46(4), 569–590. <https://doi.org/10.1080/03080188.2021.1890484>
- ◇ Bradley, T.-D., Gastaldi, J. L., & Terilla, J. (2024). The Structure of Meaning in Language: Parallel Narratives in Linear Algebra and Category Theory. *Notices of the American Mathematical Society*
- ◇ Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026, January). *Geometric Structures in \mathbb{R} -enriched adjunctions* [preprint under review]. <https://hal.science/hal-05452748>

References I

- Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623. <https://doi.org/10.1145/3442188.3445922>
- Bourdieu, P. (1979). *La distinction: Critique sociale du jugement*. Éditions de Minuit.
- Bradley, T.-D., Gastaldi, J. L., & Terilla, J. (2024). The Structure of Meaning in Language: Parallel Narratives in Linear Algebra and Category Theory. *Notices of the American Mathematical Society*.
- Chomsky, N. (1992, November). Language and the “cognitive revolutions” [Delivered November 23–27, 1992].
- Gastaldi, J. L. (2021). Why Can Computers Understand Natural Language? *Philosophy & Technology*, 34(1), 149–214. <https://doi.org/10.1007/s13347-020-00393-9>
- Gastaldi, J. L., Jarvis, S., Seiller, T., & Terilla, J. (2026, January). *Geometric Structures in \mathbb{R} -enriched adjunctions* [preprint under review]. <https://hal.science/hal-05452748>
- Gastaldi, J. L., & Pellissier, L. (2021). The Calculus of Language: Explicit Representation of Emergent Linguistic Structure through Type-Theoretical Paradigms. *Interdisciplinary Science Reviews*, 46(4), 569–590. <https://doi.org/10.1080/03080188.2021.1890484>
- Girard, J.-Y. (2011, September). *The blind spot*. European Mathematical Society.
- Hjelmslev, L. (1971a). La structure fondamentale du langage. In *Prolégomènes à une théorie du langage* [Prolégomènes à une théorie du langage] (pp. 177–231). Éditions de Minuit.
- Hjelmslev, L. (1971b). *Prolégomènes à une théorie du langage*. Éditions de Minuit.
- Kirschenbaum, M. (2023). *Again theory: A forum on language, meaning, and intent in the time of stochastic parrot*. <https://critinq.wordpress.com/2023/06/26/again-theory-a-forum-on-language-meaning-and-intent-in-the-time-of-stochastic-parrots/>

References II

- Levy, O., & Goldberg, Y. (2014). Neural word embedding as implicit matrix factorization. *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, 2177–2185.
- Marcus, S. (1967). *Introduction mathématique à la linguistique structurale*. Dunod.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *CoRR*, *abs/1310.4546*.
- Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., & Carter, S. (2020). Zoom in: An introduction to circuits [<https://distill.pub/2020/circuits/zoom-in>]. *Distill*. <https://doi.org/10.23915/distill.00024.001>
- Sennrich, R., Haddow, B., & Birch, A. (2016). Neural machine translation of rare words with subword units. *Proceedings of the 54th Annual Meeting of the ACL*, 1715–1725.
- Underwood, T. (2023, October 15). *The empirical triumph of theory* [Accessed: 2023-10-15]. <https://critiq.wordpress.com/2023/06/29/the-empirical-triumph-of-theory/>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

Université de Montréal - Université Côte d'Azur
Colloque Créativités Artificielles
Nice, France

Critique and Formalism

The Role of Formal Sciences for an Internalist Critique of LLMs

Juan Luis Gastaldi

`www.giannigastaldi.com`

ETH zürich

April 28, 2026