

Chat Token Vector
Università Ca' Foscari
Venice, Italy

Toward a Critical Formalism

Philosophical and Theoretical Effects of a Mathematical Critique of LLMs

Juan Luis Gastaldi

www.giannigastaldi.com

ETH zürich

June 12, 2025

Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Outline

Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Where Art Thou, Critique?

Where Art Thou, Critique?

- ◊ Good “**externalist**” critique

Where Art Thou, Critique?

- ◊ Good “**externalist**” critique
- ◊ Poor “**internalist**” critique

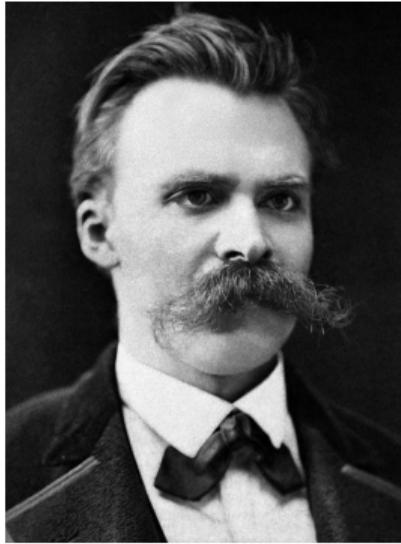
Where Art Thou, Critique?

- ◊ Good “**externalist**” critique
- ◊ Poor “**internalist**” critique
 - ◊ The main “critical” reference remains the “**Stochastic Parrots**” approach
(Bender & Koller, 2020; Bender et al., 2021)

Where Art Thou, Critique?

- ◊ Good “externalist” critique
- ◊ Poor “internalist” critique
 - ◊ The main “critical” reference remains the “**Stochastic Parrots**” approach (Bender & Koller, 2020; Bender et al., 2021)
 - ◊ **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.”

The Birth of Contemporary Critique



"In some remote corner of the universe, flickering in the light of the countless solar systems into which it had been poured, there was once a planet on which **clever animals invented cognition**. It was the most **arrogant** and most **mendacious** minute in the 'history of the world'..."

"On Truth and Lying in a Non-Moral Sense"
(Nietzsche, 1873)

The Critical Argumentative Matrix

Knowledge depends on language

The Critical Argumentative Matrix

Knowledge depends on language



The relation between language and the world is essentially arbitrary

The Critical Argumentative Matrix

Knowledge depends on language



The relation between language and the world is essentially arbitrary



Any regularity in language/knowledge is not natural but cultural/social/political

The Critical Argumentative Matrix

Knowledge depends on language



The relation between language and the world is essentially arbitrary



Any regularity in language/knowledge is not natural but cultural/social/political



We should resist existing regularities and create new ones

The Critical Argumentative Matrix

Knowledge depends on language
(Epistemological)



The relation between language and the world is essentially arbitrary



Any regularity in language/knowledge is not natural but cultural/social/political
(Political)



We should resist existing regularities and create new ones
(Aesthetic)

The Critical Argumentative Matrix

Knowledge depends on language
(Epistemological)

[The relation between language and the world is essentially arbitrary?]

Any regularity in language/knowledge is not natural but cultural/social/political
(Political)

We should resist existing regularities and create new ones
(Aesthetic)

- ◊ At the source of this situation is the new foundational role played by **formal sciences** in the 20th century
 - ◊ For a **theory of language**: Carnap, Gödel, Turing, Shannon, Harris, Chomsky...

- ◊ At the source of this situation is the new foundational role played by **formal sciences** in the 20th century
 - ◊ For a **theory of language**: Carnap, Gödel, Turing, Shannon, Harris, Chomsky...
- ◊ The critical tradition has either **withdrawn** from the areas conquered by formal approaches, or made formal approaches the **target** of criticism

- ◊ At the source of this situation is the new foundational role played by **formal sciences** in the 20th century
 - ◊ For a **theory of language**: Carnap, Gödel, Turing, Shannon, Harris, Chomsky...
- ◊ The critical tradition has either **withdrawn** from the areas conquered by formal approaches, or made formal approaches the **target** of criticism
- ◊ We need a **new strategy**: Elaborate a **critical formalism**

- ◊ In the case of **AI**, a critical formalism can provide:
 - ◊ New **epistemological tools** countering dogmatic perspectives stemming from within the field
 - ◊ New **theoretical tools** contributing to the non-dogmatic positive production of knowledge

Outline

Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Neural LMs as Computable Functions

Neural LM



?

Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions

Neural LM



?

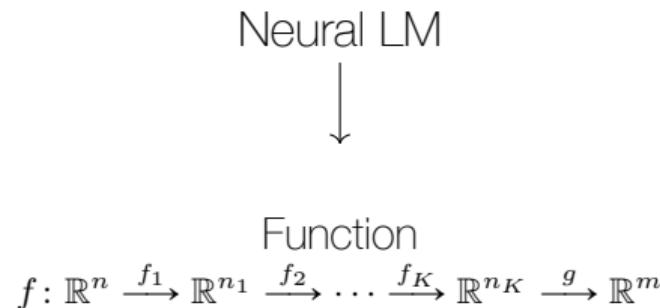
Neural LMs as Computable Functions

Neural LM

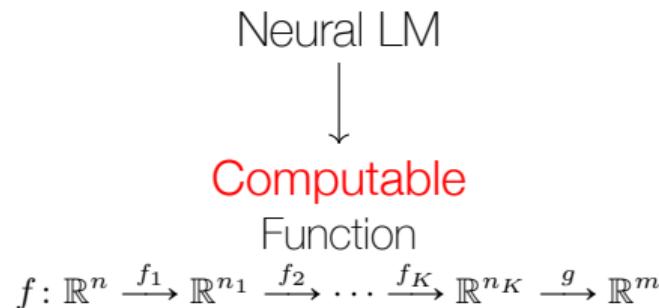


f !

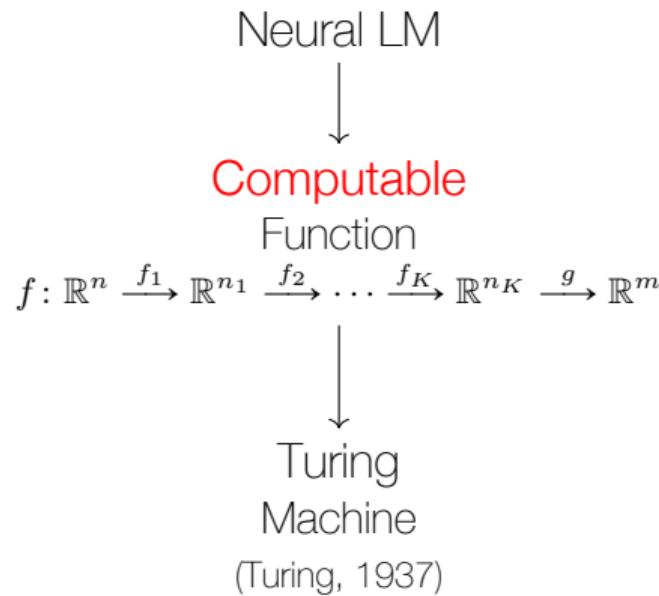
Neural LMs as Computable Functions



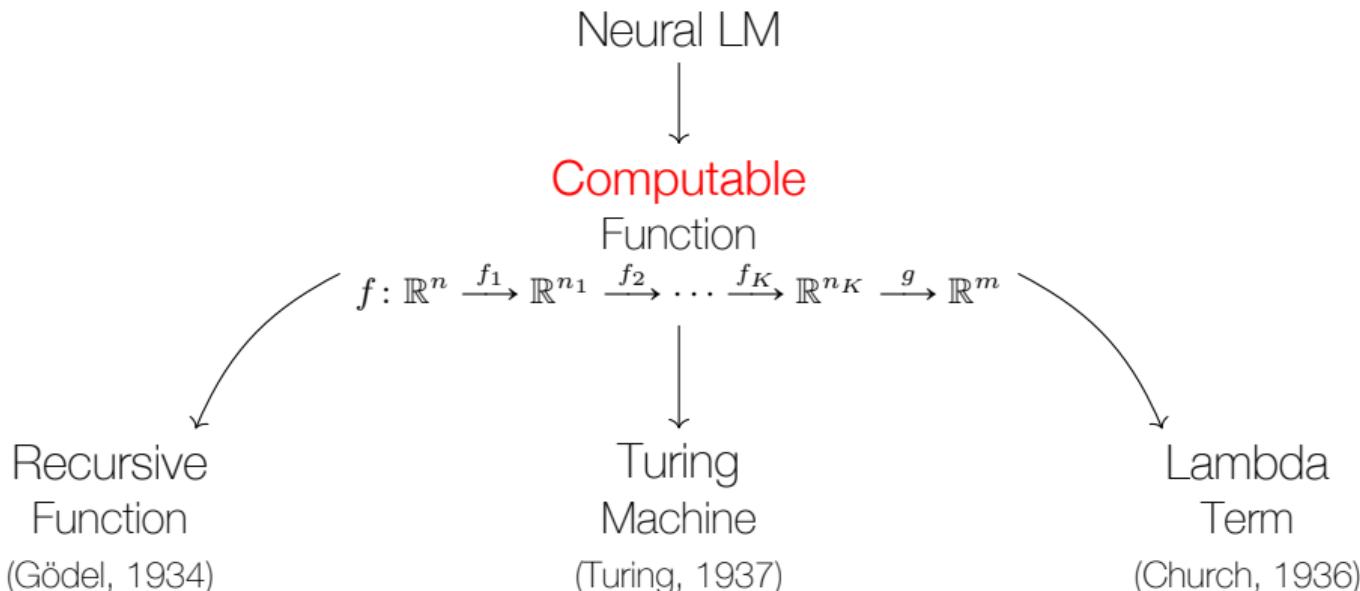
Neural LMs as Computable Functions



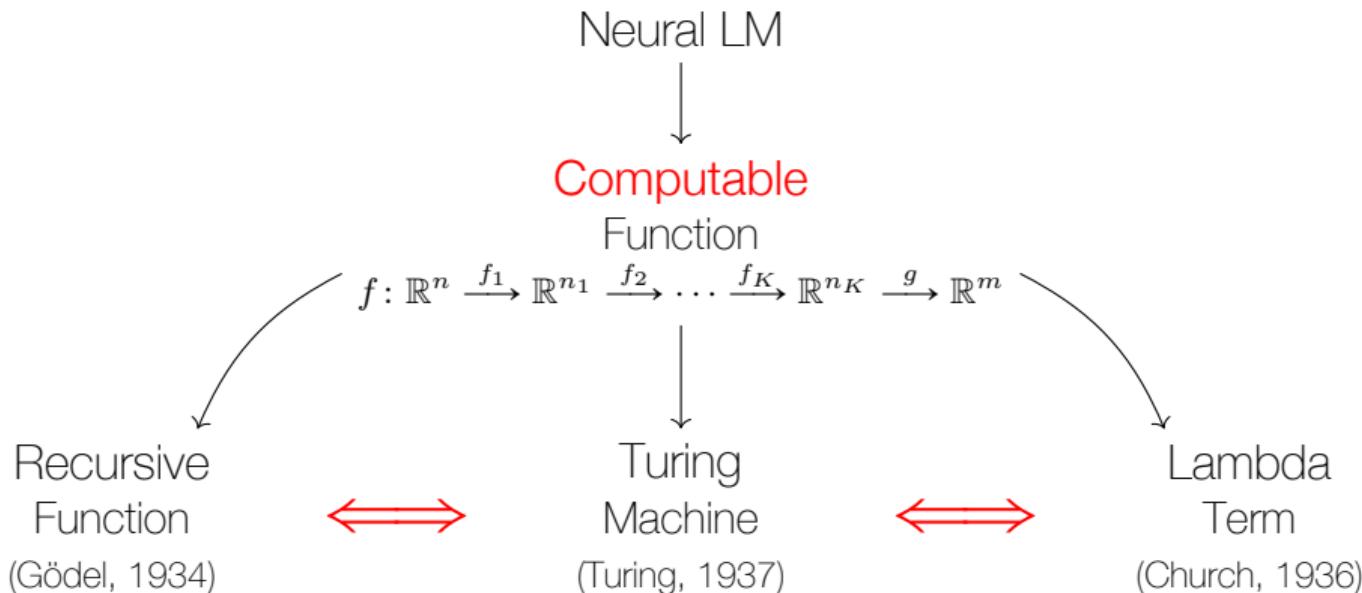
Neural LMs as Computable Functions



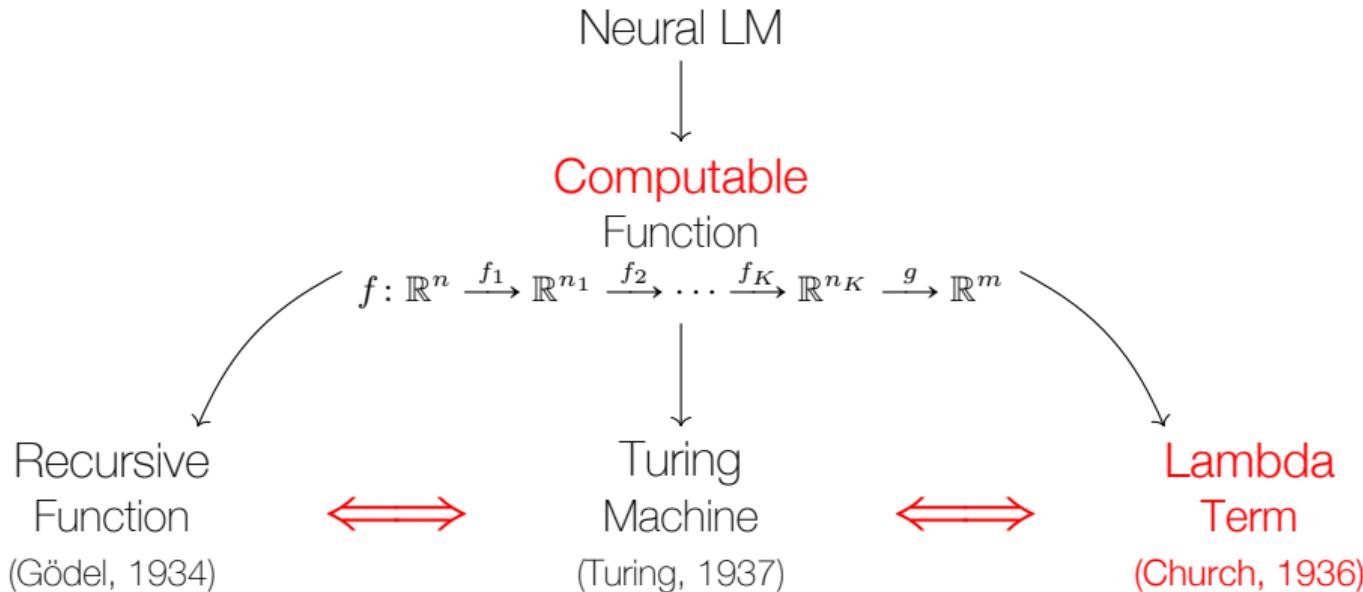
Neural LMs as Computable Functions



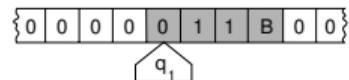
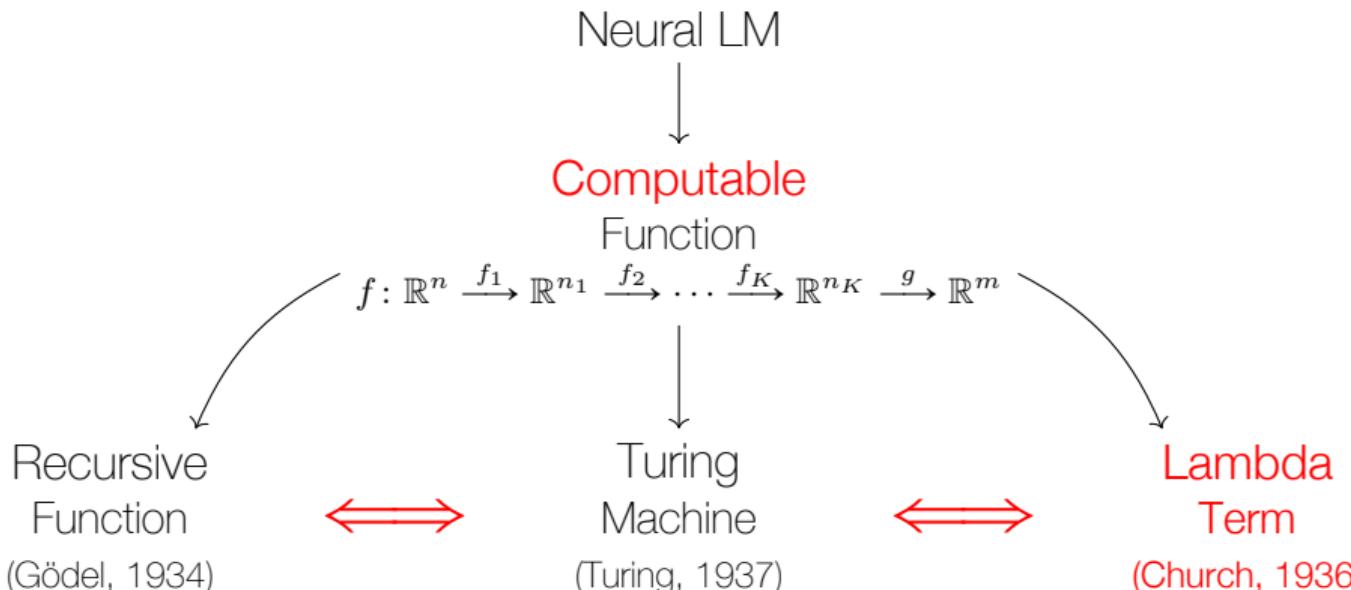
Neural LMs as Computable Functions



Neural LMs as Computable Functions



Neural LMs as Computable Functions



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

credit: Nynexman4464

λ -abstraction and β -reduction in λ -calculus

yxz

λ -abstraction and β -reduction in λ -calculus

$$\lambda \color{red}x.y\color{black}xz$$

λ -abstraction and β -reduction in λ -calculus

$$(\lambda \textcolor{red}{x}.y \textcolor{red}{x} z) \textcolor{blue}{t}$$

λ -abstraction and β -reduction in λ -calculus

$$(\lambda \textcolor{red}{x}.y \textcolor{red}{x} z) \textcolor{blue}{t}$$

$$y \textcolor{blue}{t} z$$

Empirical Evaluation

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

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

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

...

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

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

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

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

$$P' := \color{blue}{\lambda r. \lambda s. \lambda f. \lambda x. f(f(f(f(fx))))}$$

0:	$\lambda f. \lambda x. x$	$\color{blue}{\lambda r. \lambda s. \lambda f. \lambda x. f(f(f(f(fx))))} (\color{orange}{\lambda f. \lambda x. f(fx)}) (\color{green}{\lambda f. \lambda x. f(f(fx))})$
1:	$\lambda f. \lambda x. f x$	↓
2:	$\color{orange}{\lambda f. \lambda x. f(fx)}$	↓
3:	$\color{green}{\lambda f. \lambda x. f(f(fx))}$	↓
4:	$\lambda f. \lambda x. f(f(f(fx)))$	↓
5:	$\color{red}{\lambda f. \lambda x. f(f(f(f(fx))))}$	↓
...		↓
n:	$\lambda f. \lambda x. \underbrace{f(\dots(f}_{n \text{ times}} x) \dots)$	$\color{red}{\lambda f. \lambda x. f(f(f(f(f(fx)))))}$

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

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

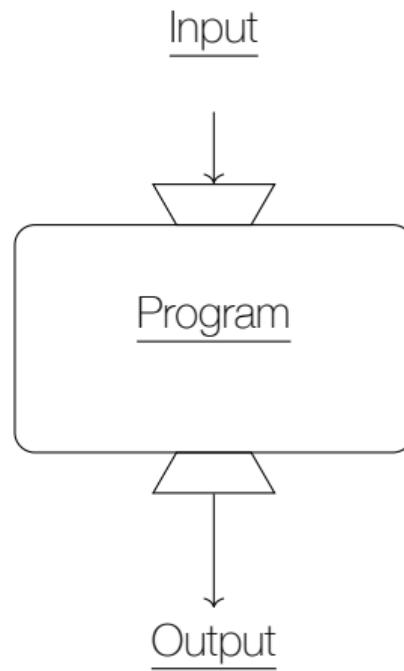
$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f(n f x)$$

0:	$\lambda f. \lambda x. x$	$\lambda m. \lambda n. \lambda f. \lambda x. m f(n f x) (\lambda f. \lambda x. f(fx)) (\lambda f. \lambda x. f(f(fx)))$
1:	$\lambda f. \lambda x. f x$	$\lambda m. \lambda n. \lambda f. \lambda x. m f(n f x) (\lambda g. \lambda y. g(gy)) (\lambda h. \lambda z. h(h(hz)))$
2:	$\lambda f. \lambda x. f(fx)$	$\lambda n. \lambda f. \lambda x. (\lambda g. \lambda y. g(gy)) f(n f x) (\lambda h. \lambda z. h(h(hz)))$
3:	$\lambda f. \lambda x. f(f(fx))$	$\lambda n. \lambda f. \lambda x. (\lambda g. \lambda y. g(gy)) f(n f x) (\lambda h. \lambda z. h(h(hz)))$
4:	$\lambda f. \lambda x. f(f(f(fx))))$	$\lambda f. \lambda x. (\lambda g. \lambda y. g(gy)) f((\lambda h. \lambda z. h(h(hz))) f x)$
5:	$\lambda f. \lambda x. f(f(f(f(fx))))$	$\lambda f. \lambda x. (\lambda y. f(fy)) ((\lambda h. \lambda z. h(h(hz))) f x)$
...		$\lambda f. \lambda x. (\lambda y. f(fy)) ((\lambda z. f(f(fz))) x)$
$n:$	$\lambda f. \lambda x. \underbrace{f(\dots(f}_{n \text{ times}} x) \dots)$	$\lambda f. \lambda x. (\lambda y. f(fy)) (f(f(fx)))$
		$\lambda f. \lambda x. f(f(f(f(fx)))))$

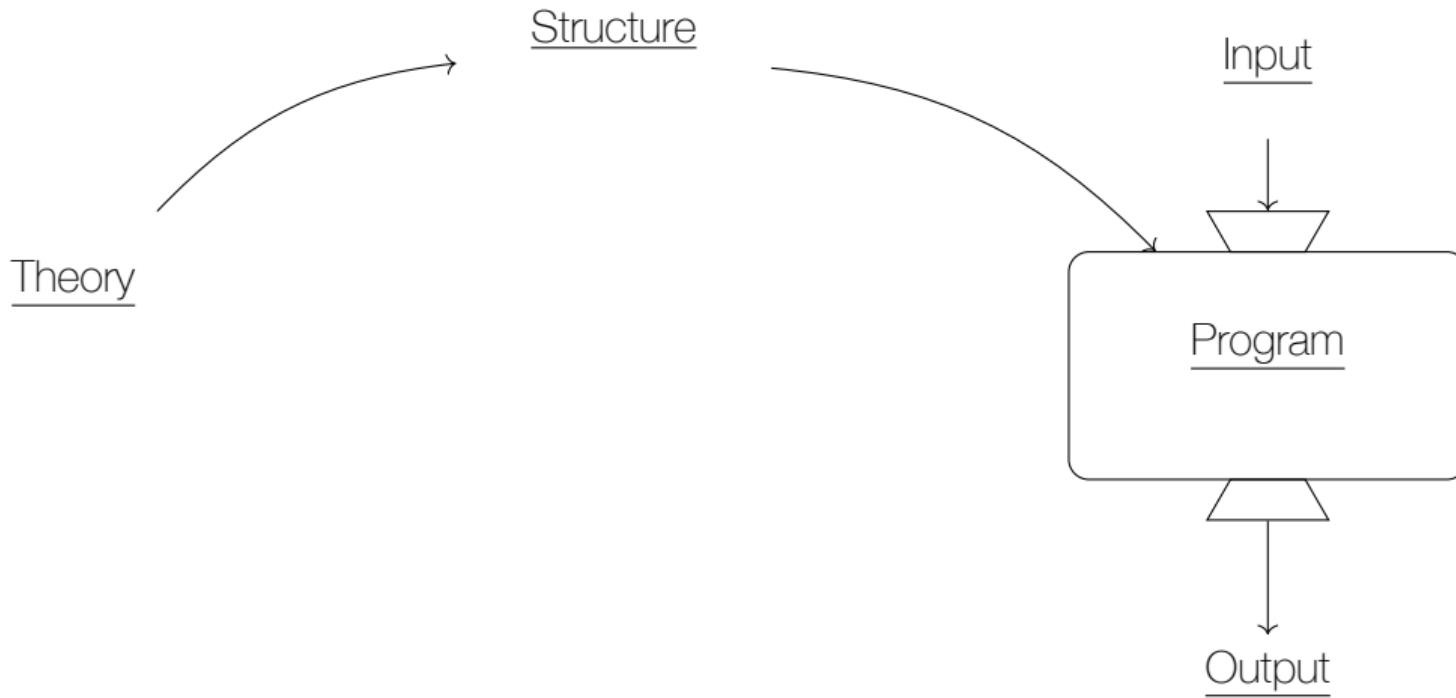
$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$$

$P'' := \lambda R o f \tilde{A} O e \tilde{N} 5 \tilde{E} | \tilde{A} x \tilde{x} = \infty \tilde{u} \tilde{y} m W f 286 \tilde{e} y' S \tilde{O} \tilde{u} > v \& \tilde{i} \tilde{A} \neg 2 \tilde{o} \tilde{E} 7 \tilde{o} \tilde{c} \tilde{\infty} \{ \tilde{a} > 2 \tilde{f} \tilde{l} \tilde{B} \tilde{u} \tilde{G} \# \tilde{A} 9 \tilde{C} \tilde{U}$
 $\infty \tilde{b} t \tilde{Y} \tilde{B} \tilde{b} \tilde{Y} \tilde{U} \tilde{e} \% 3 ; 5 \tilde{a} [\tilde{l} - \tilde{e} u \tilde{o} \tilde{U} \tilde{.} \tilde{7} - \tilde{U} . \lambda : \tilde{4} \tilde{m} \tilde{O} \tilde{O} \tilde{Y} \tilde{e} \tilde{-} + \tilde{I} \tilde{s} \tilde{O} \tilde{,} \tilde{S} \tilde{+} \tilde{g} \tilde{i} \tilde{,} \tilde{B} \tilde{T} \tilde{M} \tilde{\div} \tilde{o} \tilde{-} \# \tilde{i} \tilde{Y} \tilde{e} \tilde{U} \tilde{v}$
 $- g \tilde{O} \tilde{y} / \tilde{e} i i j \tilde{O} \tilde{t} \tilde{C} \tilde{f} \tilde{i} \tilde{f} \tilde{.} \tilde{J} \tilde{1} \tilde{«} \tilde{e} \tilde{\emptyset} \tilde{,} \tilde{I} \tilde{h} \tilde{a} \tilde{e} \tilde{t} \tilde{f} \tilde{a} \tilde{e} \tilde{Y} \tilde{S} \tilde{^} \tilde{6} \tilde{F} \tilde{i} \tilde{W} \tilde{»} \tilde{R} \tilde{U} \tilde{K} \tilde{g} \tilde{e} \tilde{.} \tilde{\lambda} \tilde{f} \tilde{d} \tilde{-} \dots \tilde{D} \tilde{2} \tilde{\div} \tilde{o} \tilde{.} \tilde{x} \tilde{e} \tilde{E} \tilde{y} \tilde{.} \tilde{O} \tilde{”} \tilde{c} \tilde{b}$
 $B \tilde{e} \tilde{f} N \tilde{E} 1 \tilde{E} \tilde{f} / \tilde{U} \tilde{9} \tilde{N} \tilde{p} \tilde{u} / \tilde{J} \tilde{Y} \tilde{C} \tilde{o} \tilde{E} 9 \tilde{y} \tilde{A} \tilde{E} \tilde{.} \tilde{\lambda} \tilde{A} \tilde{I} \tilde{A} \tilde{^} \tilde{o} \tilde{C} \tilde{,} \tilde{»} \tilde{f} \tilde{q} \tilde{\infty} \tilde{\pm} \tilde{i} \tilde{B} \tilde{5} \tilde{l} \tilde{>} \tilde{O} \tilde{”} \tilde{g} \tilde{T} \tilde{M} \tilde{“} \tilde{6} \tilde{\Omega} \tilde{e} \tilde{“} \tilde{a} \tilde{e} \tilde{e} \tilde{C} \tilde{/} \tilde{a} \tilde{...} \tilde{O} \tilde{.} \tilde{f} \tilde{O} \tilde{A} \tilde{] \tilde{N} \tilde{a} \tilde{y} \tilde{E} \tilde{N} \tilde{^} \tilde{E} \tilde{.} \tilde{»} \tilde{(} \tilde{f} \tilde{d} \tilde{-} \dots \tilde{D} \tilde{2} \tilde{\div} \tilde{o} \tilde{.} \tilde{x} \tilde{e} \tilde{E} \tilde{y} \tilde{.} \tilde{O} \tilde{”} \tilde{c} \tilde{b} B \tilde{e} \tilde{f} N \tilde{E} 1 \tilde{E} \tilde{f} / \tilde{U} \tilde{9} \tilde{N} \tilde{p} \tilde{u} / \tilde{J} \tilde{Y} \tilde{C} \tilde{o} \tilde{E} 9 \tilde{y} \tilde{A} \tilde{E} \tilde{A} \tilde{I} \tilde{A} \tilde{^} \tilde{o} \tilde{C} \tilde{,} \tilde{»} \tilde{f} \tilde{q} \tilde{\infty} \tilde{\pm} \tilde{i} \tilde{B} \tilde{5} \tilde{l} \tilde{>} \tilde{O} \tilde{”} \tilde{g} \tilde{T} \tilde{M} \tilde{“} \tilde{6} \tilde{\Omega} \tilde{e} \tilde{“} \tilde{a} \tilde{e} \tilde{e} \tilde{C} \tilde{/} \tilde{a} \tilde{...} \tilde{O} \tilde{.} \tilde{f} \tilde{O} \tilde{A} \tilde{] \tilde{N} \tilde{a} \tilde{y} \tilde{E} \tilde{N} \tilde{^} \tilde{E} \tilde{.} \tilde{»} \tilde{A} \tilde{à} \tilde{e} \tilde{f} U \tilde{ò} \tilde{f} E \tilde{U} \tilde{.} \tilde{I} \tilde{m} \tilde{\#} \tilde{,} \tilde{,} \tilde{4} \tilde{\backslash} \tilde{r} \tilde{\sqrt{}} \tilde{-} \tilde{\div} \tilde{\tilde{I} \tilde{p} \tilde{o}} \tilde{»} \tilde{y} \tilde{*} \tilde{v} \tilde{t} \tilde{\tilde{A} \tilde{J} \tilde{A} \tilde{F} \tilde{1} \tilde{u} \tilde{A} \tilde{ó} \tilde{z} \tilde{«} \tilde{ñ} \tilde{M} \tilde{”} \tilde{D} \tilde{j} \tilde{C} \tilde{E} B \tilde{E} \tilde{è} \tilde{Í} \tilde{T} \tilde{—} \tilde{E} \tilde{a} \tilde{\%} \tilde{A} \tilde{C} \tilde{\Omega} \tilde{@} \tilde{\backslash} \tilde{\backslash} \tilde{O} \tilde{\wedge} \tilde{~} \tilde{]} \tilde{\tilde{I} \tilde{h} \tilde{f}} \tilde{)} \tilde{(} \tilde{\tilde{E} \tilde{I} \tilde{U} \tilde{e} \tilde{í} \tilde{4} \tilde{W} \tilde{p} \tilde{í}} \tilde{\}} \tilde{w} \tilde{,} \tilde{\$} \tilde{\Omega} \tilde{“} \tilde{K} \tilde{5} \tilde{e} \tilde{A} \tilde{\P} \tilde{\%} \tilde{3} \tilde{[} \tilde{m} \tilde{,} \tilde{”} \tilde{B} \tilde{A} \tilde{f} \tilde{f} \tilde{O} \tilde{;} \tilde{o} \tilde{J} \tilde{ç} \tilde{C} \tilde{E} \tilde{í} \tilde{o} \tilde{Y} \tilde{O} \tilde{c} \tilde{B} \tilde{,} \tilde{\$} \tilde{A} \tilde{á} \tilde{\}} \tilde{O} \tilde{A} \tilde{\%} \tilde{3} \tilde{;}$
 $\tilde{?} \tilde{o} \tilde{-} \tilde{o} \tilde{C} \tilde{E} \tilde{@} \tilde{f} \tilde{l} \tilde{8} \tilde{”} \tilde{R} \tilde{C} \tilde{æ} \tilde{e} \tilde{o} \tilde{*} \tilde{&} \tilde{<} \tilde{Y} \tilde{-} \tilde{o} \tilde{1} \tilde{2} \tilde{A} \tilde{\%} \tilde{a} \tilde{O} \tilde{Ü} \tilde{\#} \tilde{i} \tilde{”} \tilde{,} \tilde{ú} \tilde{”} \tilde{«} \tilde{\hat{o}} \tilde{,} \tilde{\infty} \tilde{I} \tilde{a} \tilde{ä} \tilde{“} \tilde{\phi} \tilde{A} \tilde{d} \tilde{|} \tilde{”} \tilde{N} \tilde{’} \tilde{E} \tilde{y} \tilde{\ø} \tilde{;} \tilde{”} \tilde{W} \tilde{»} \tilde{w} \tilde{o} \tilde{[} \tilde{]} \tilde{\»} \tilde{\tilde{O} \tilde{E} \tilde{u} \tilde{w} \tilde{’} \tilde{6} \tilde{<} \tilde{ù} \tilde{”} \tilde{=} \tilde{\tilde{a} \tilde{O} \tilde{-} \tilde{I} \tilde{D} \tilde{z} \tilde{?} \tilde{2} \tilde{\pm} \tilde{|} \tilde{é} \tilde{’} \tilde{3} \tilde{A} \tilde{/} \tilde{r} \tilde{x} \tilde{\mu} \tilde{\infty} \tilde{\mu} \tilde{\$} \tilde{\tilde{A} \tilde{e} \tilde{A} \tilde{*} \tilde{f} \tilde{l} \tilde{”} \tilde{\hat{u}} \tilde{’} \tilde{+} \tilde{I} \tilde{V} \tilde{y} \tilde{a} \tilde{G} \tilde{æ} \tilde{ß} \tilde{ä} \tilde{g} \tilde{\hat{o}} \tilde{/} \tilde{,} \tilde{u} \tilde{N}}$

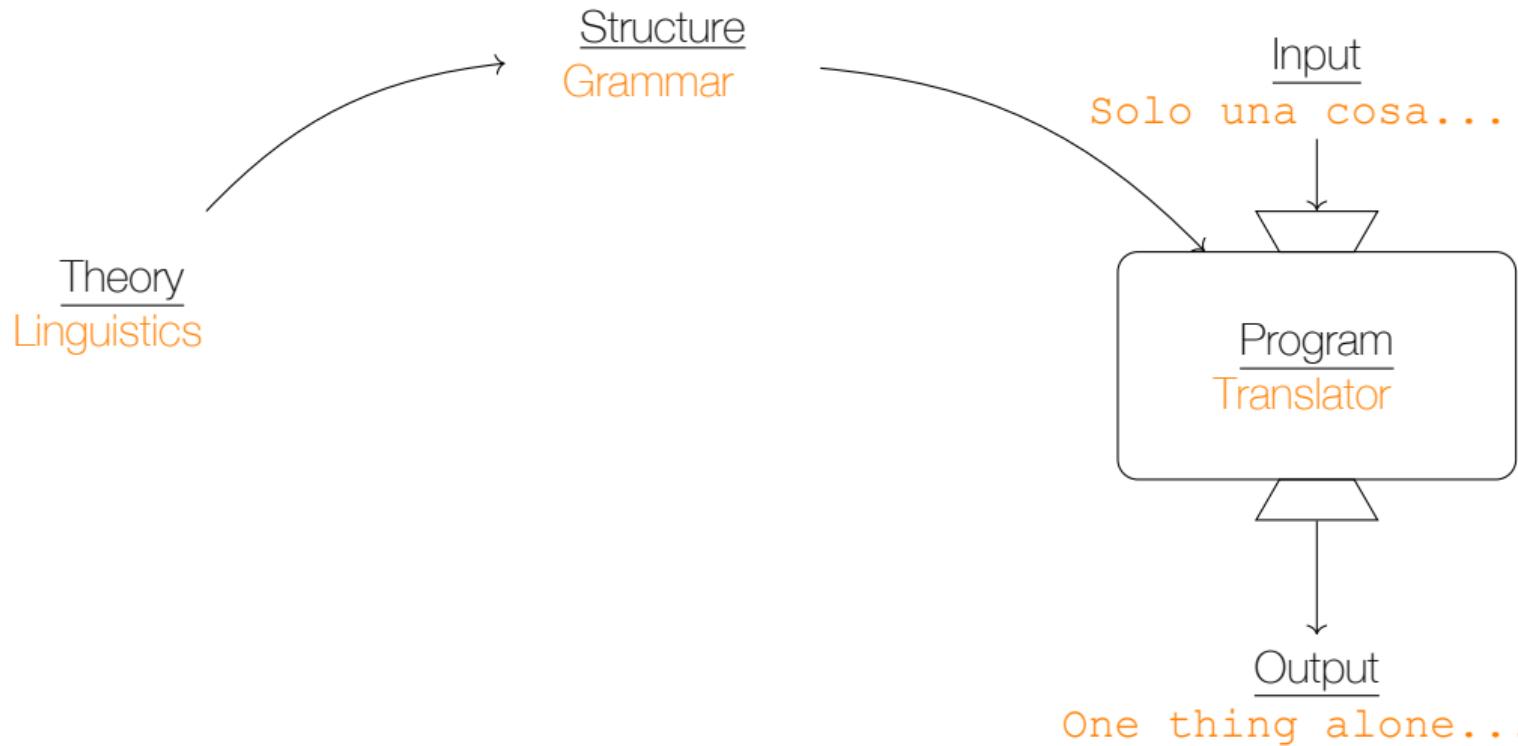
Making It Explicit



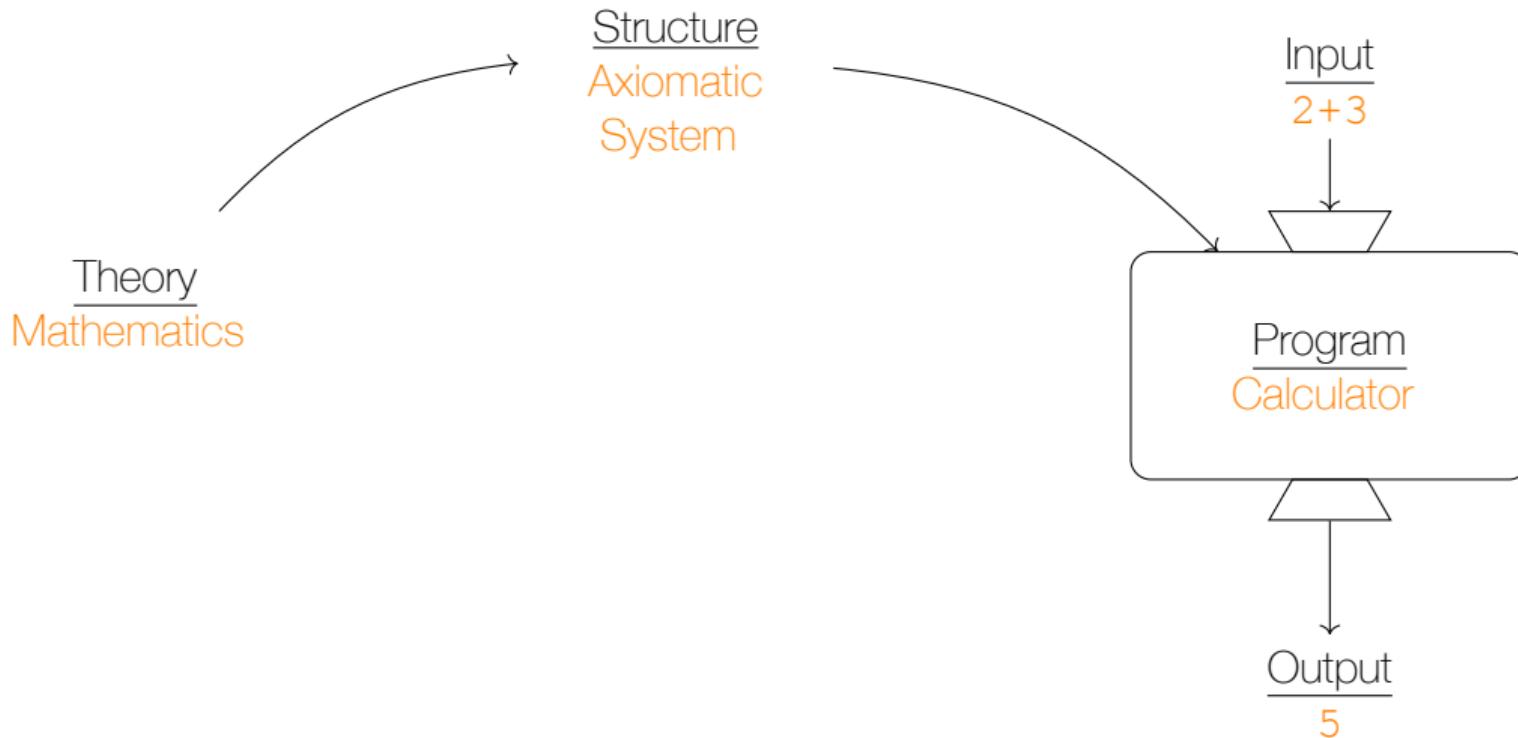
Making It Explicit



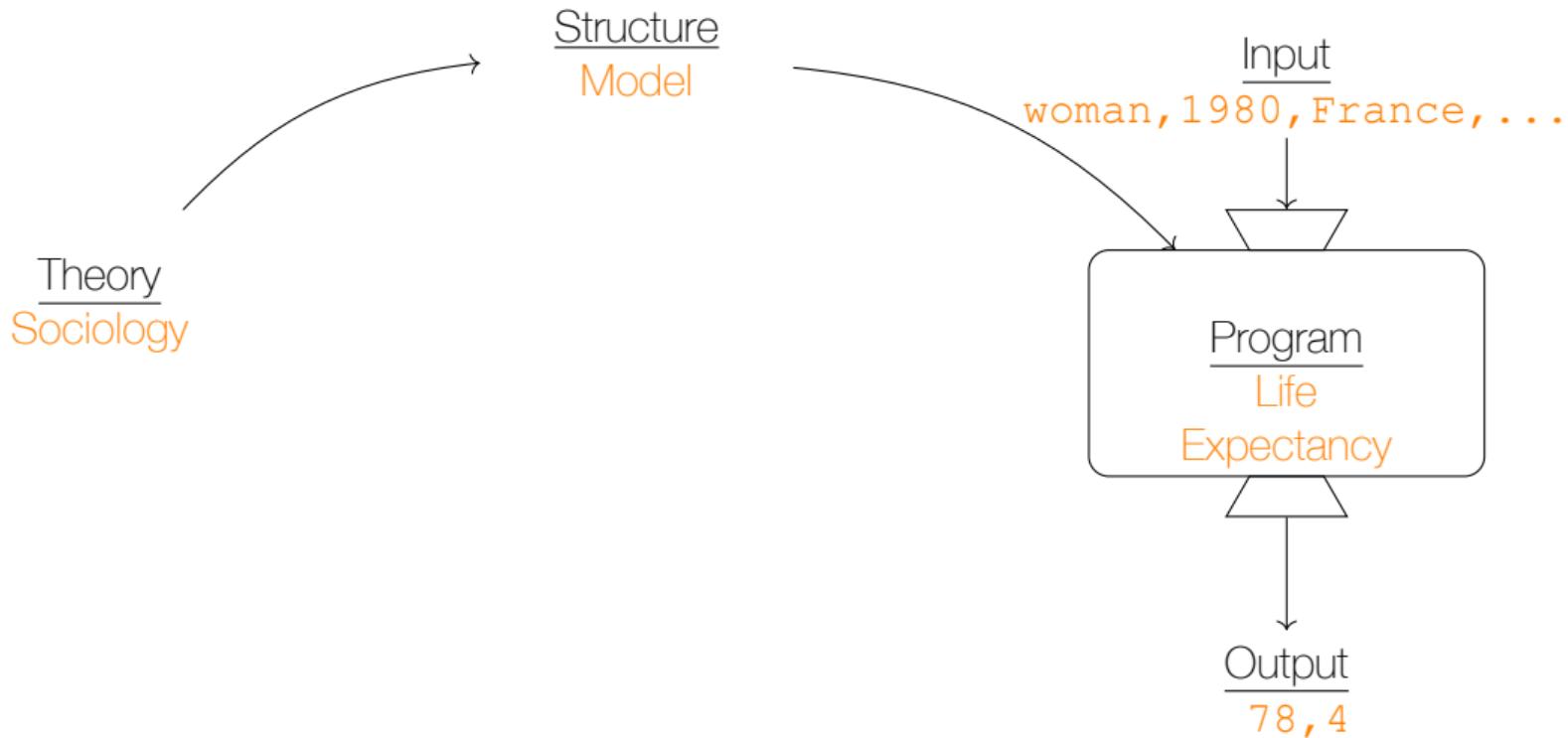
Making It Explicit



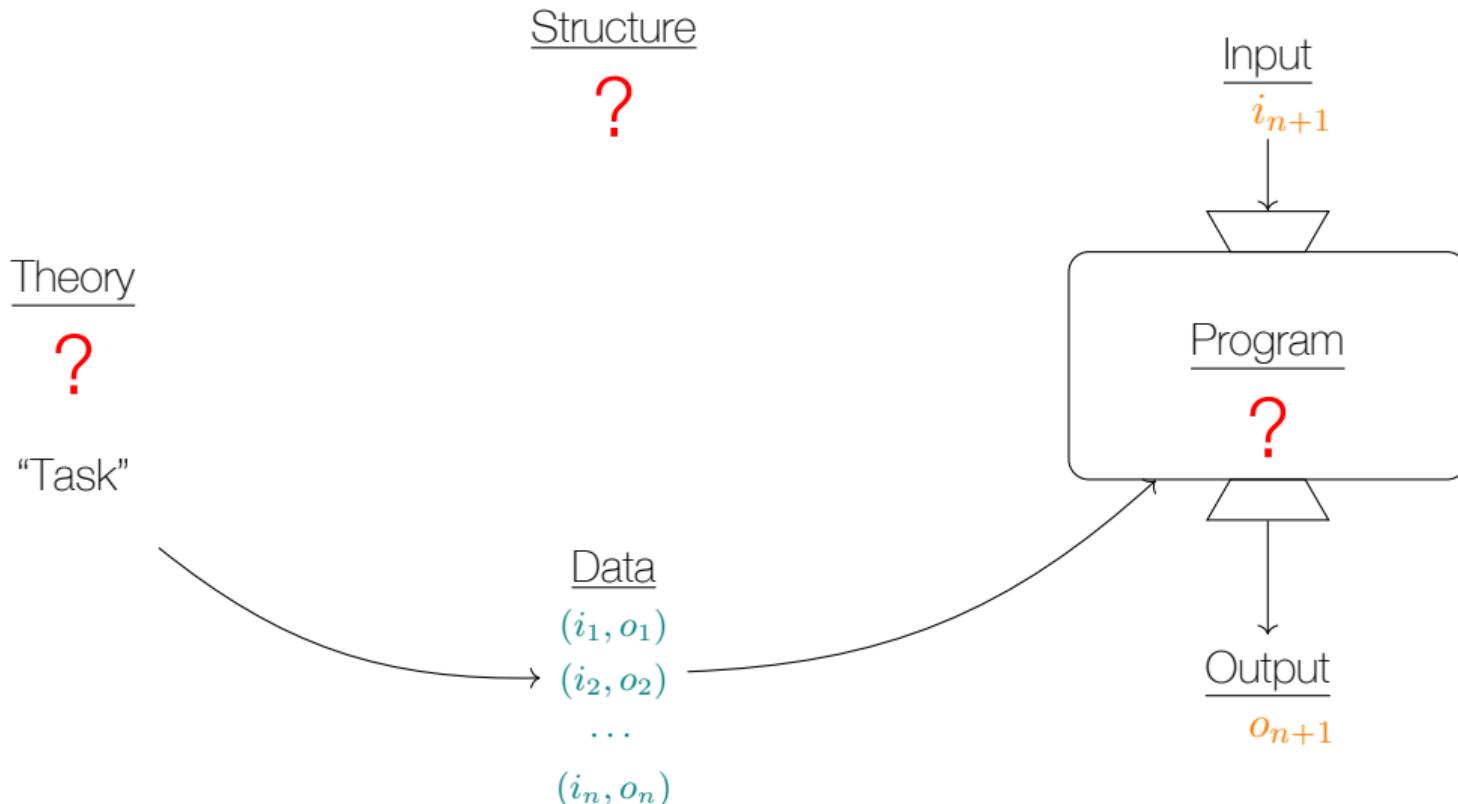
Making It Explicit



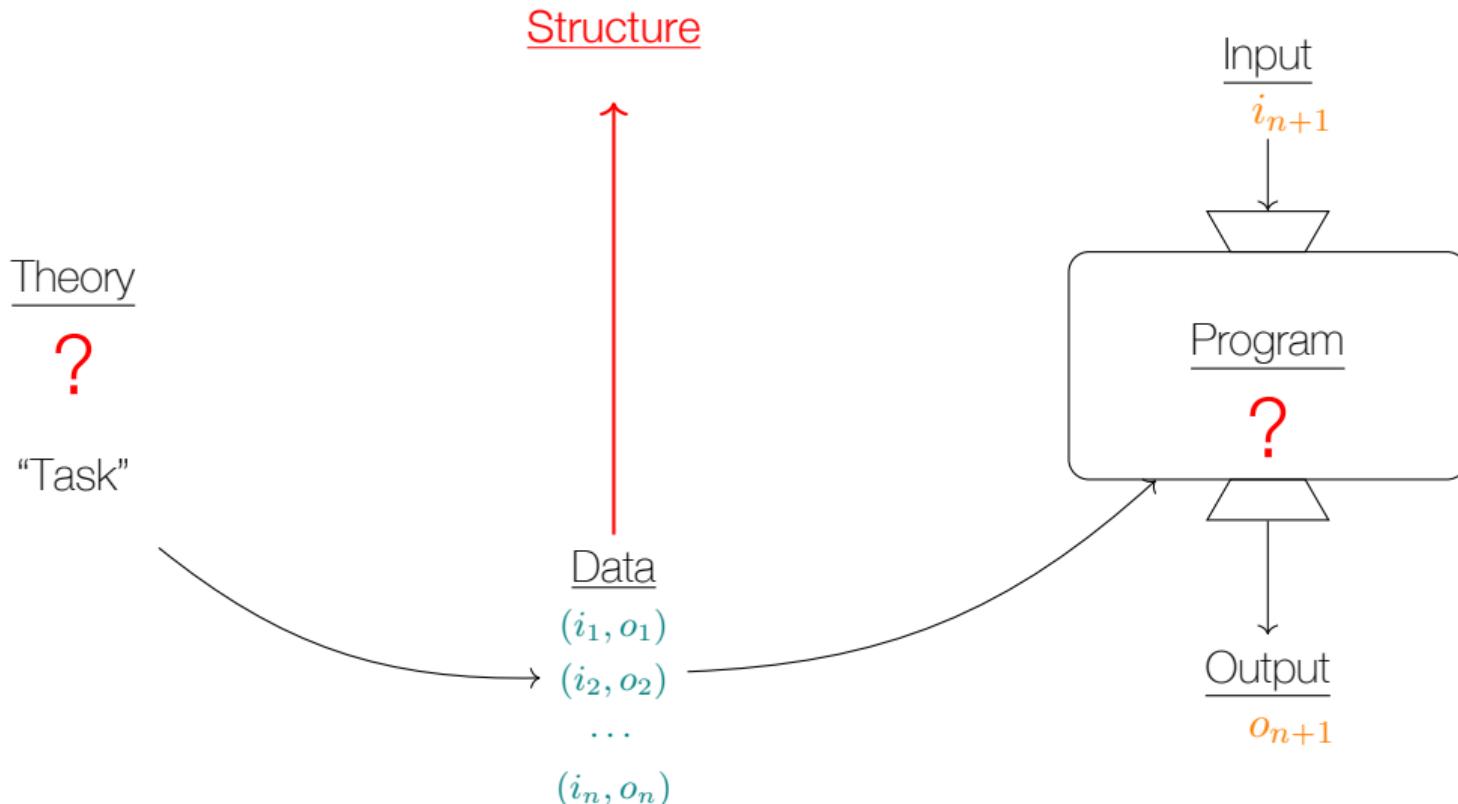
Making It Explicit



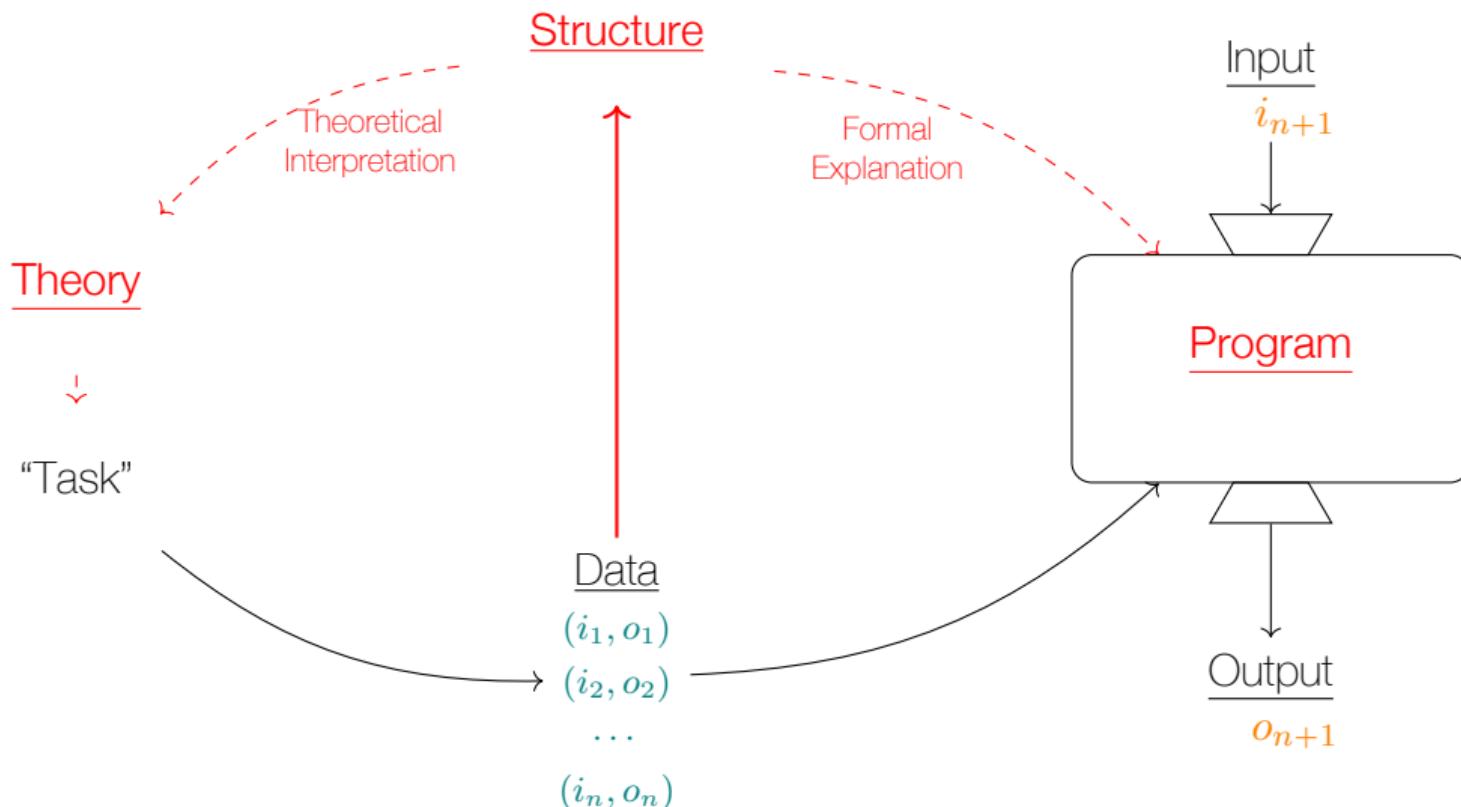
Making It Explicit



Making It Explicit



Making It Explicit



Outline

Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

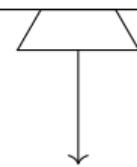
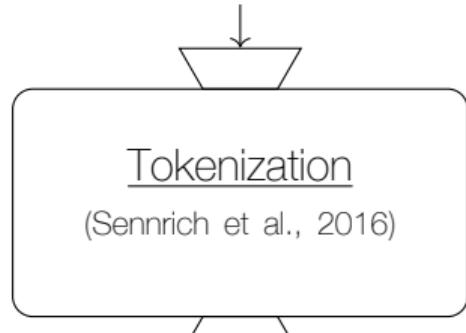
The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Formal Explainability

Epistemology of Machine Learning
Distributional Language Models

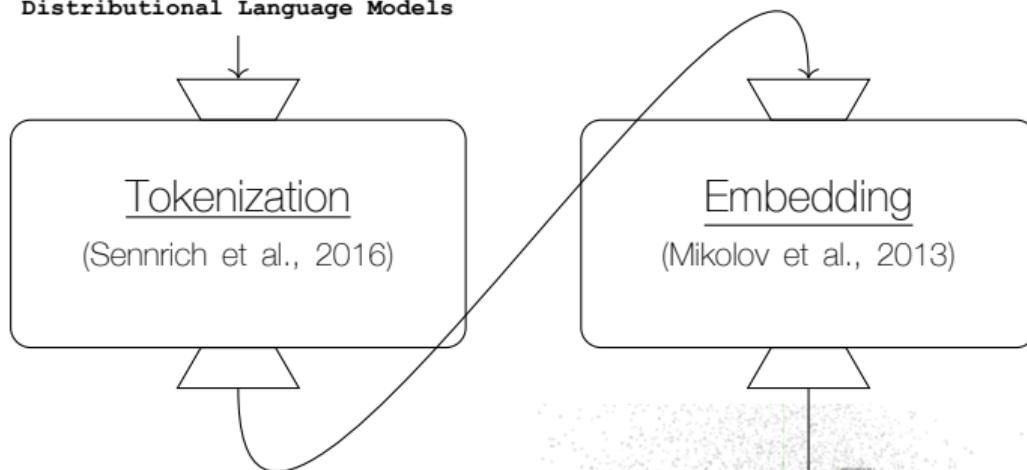


Epistemology of Machine Learning
Distributional Language Models

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

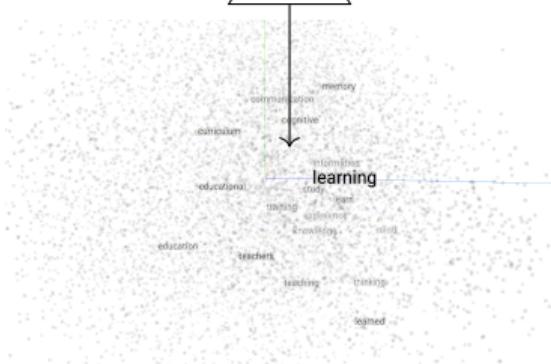
Formal Explainability

Epistemology of Machine Learning Distributional Language Models



Epistemology of Machine Learning Distributional Language Models

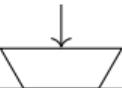
(<https://tiktok-encoder.vercel.app>)



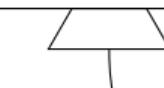
(<https://projector.tensorflow.org>)

Formal Explainability

**Epistemology of Machine Learning
Distributional Language Models**



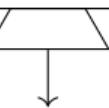
Tokenization
(Sennrich et al., 2016)



Embedding
(Mikolov et al., 2013)

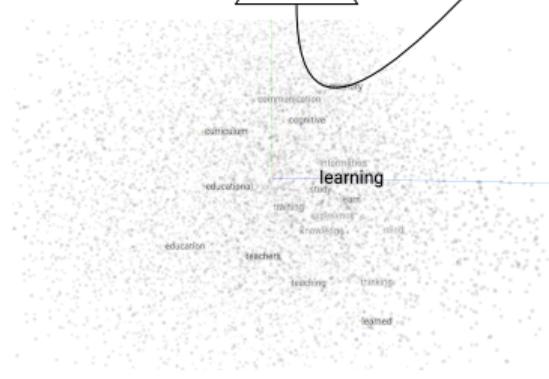


Attention
(Vaswani et al., 2017)



**Epistemology of Machine Learning
Distributional Language Models**

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



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

Ep
ist
em
olog
y
of
Ma
chi
ne
Le
ar
ni
ng

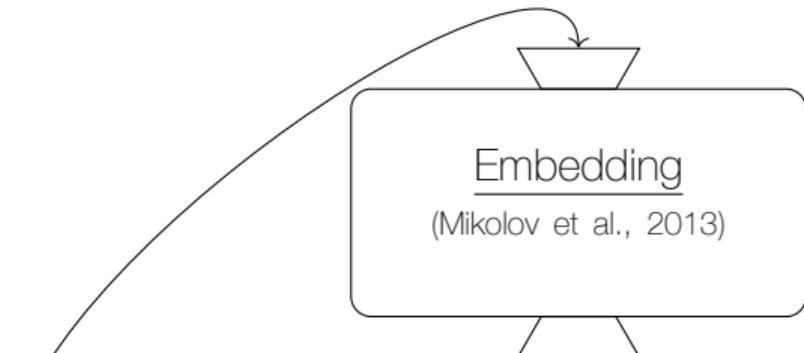
Distr
ib
ut
io
n
al
La
ng
u
ag
e
Mo
de
ls

Ep
ist
em
olog
y
of
Ma
chi
ne
Le
ar
ni
ng

Distr
ib
ut
io
n
al
La
ng
u
ag
e
Mo
de
ls

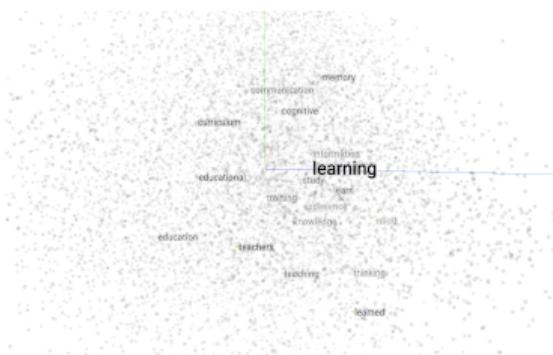
(<https://github.com/jessevieg/bertviz>)

Formal Explainability



**Epistemology of Machine Learning
Distributional Language Models**

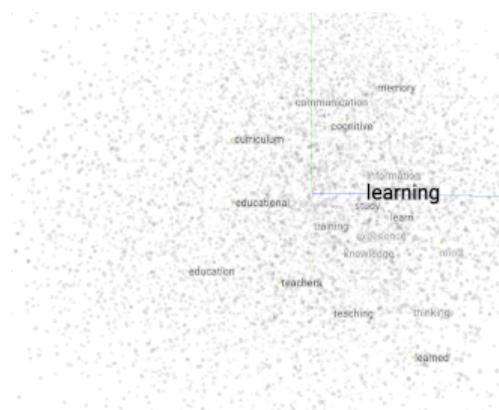
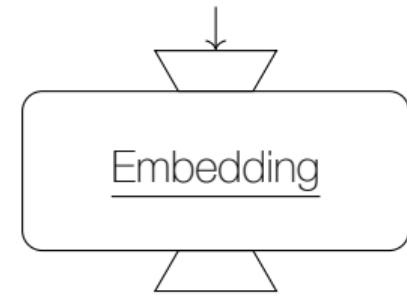
(<https://tiktoktokenizer.vercel.app>)



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

The Structure of Embeddings

Epistemology of Machine Learning
Distributional Language Models

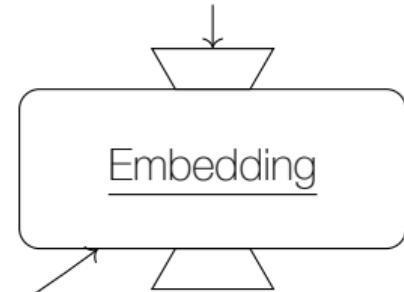
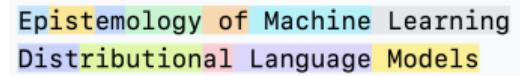
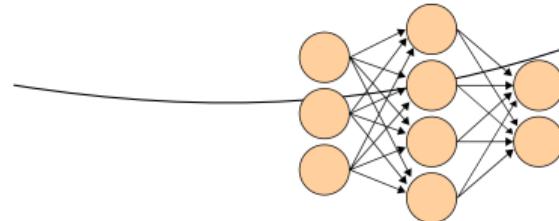


The Structure of Embeddings

Structure

?

Data

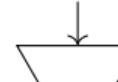


The Structure of Embeddings

Structure

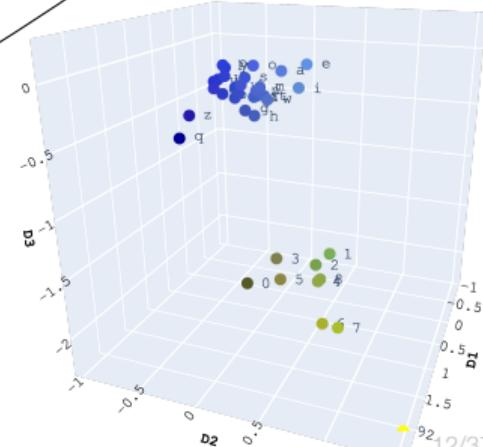
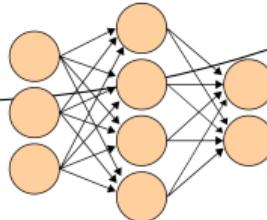
?

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



Embedding

Data



Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an implicit, low-dimensional factorization of a pointwise mutual information (pmi), word-context matrix.

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an **implicit**, low-dimensional factorization of a pointwise mutual information (pmi), word-context matrix.

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an **implicit, low-dimensional** factorization of a pointwise mutual information (pmi), word-context matrix.

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an **implicit, low-dimensional factorization** of a pointwise mutual information (pmi), word-context matrix.

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an **implicit, low-dimensional factorization** of a **pointwise mutual information (pmi)**, word-context matrix.

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an **implicit, low-dimensional factorization** of a **pointwise mutual information (pmi), word-context** matrix.

word2vec Explained (Levy and Goldberg, 2014)

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

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

- Word2vec performs an **implicit, low-dimensional factorization** of a **pointwise mutual information (pmi), word-context matrix.**

word2vec Explained (Levy and Goldberg, 2014)

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

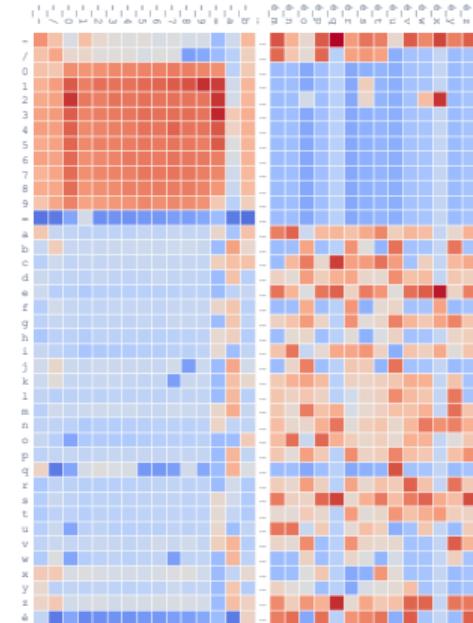
$$\frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{when} \quad \vec{w} \cdot \vec{c} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

- ◊ 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.

Example: Characters in Wikipedia

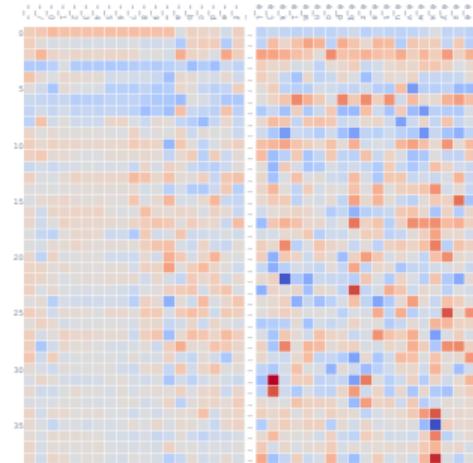
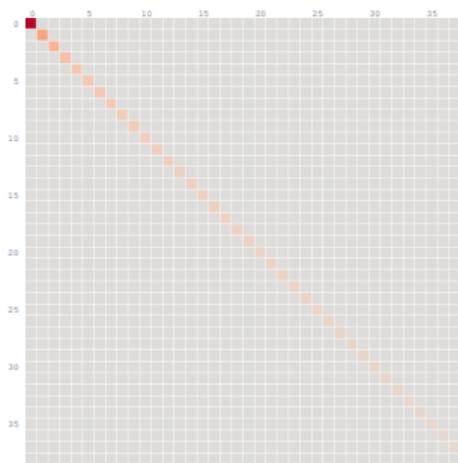
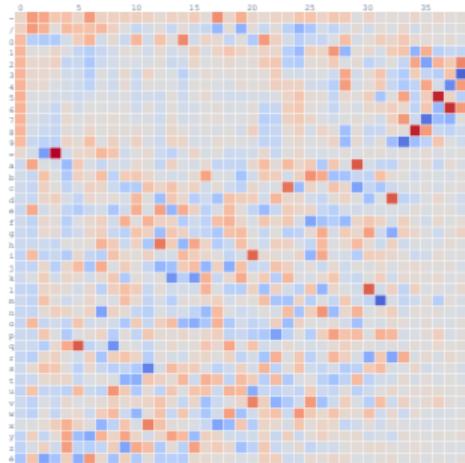
$$W = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, é\}$$

$$C = X \times X = \{ (-, -), (-, /), (-, 0), \dots, (é, z), (é, é) \}$$



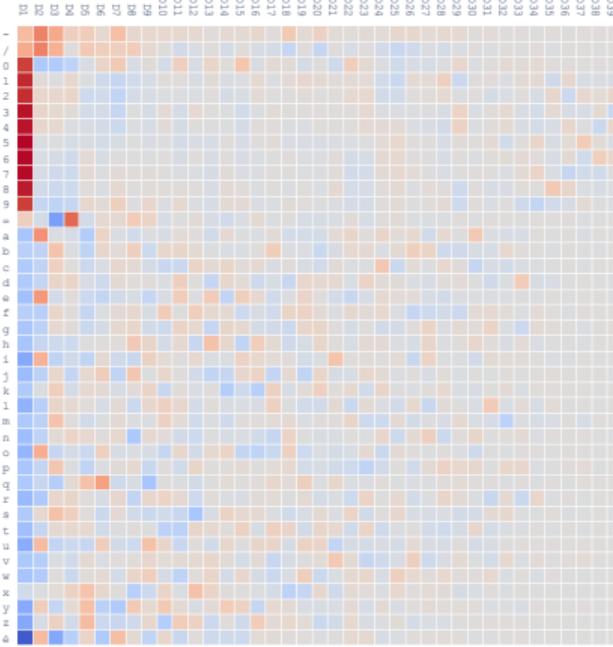
$$\begin{aligned} M_{wc} &= \text{pmi}(w, c) \\ &= \log \frac{p(w, c)}{p(w)p(c)} \end{aligned}$$

SVD of Wikipedia Character PMI Matrix

 U Σ V^T 

Truncate

$U \times \Sigma$



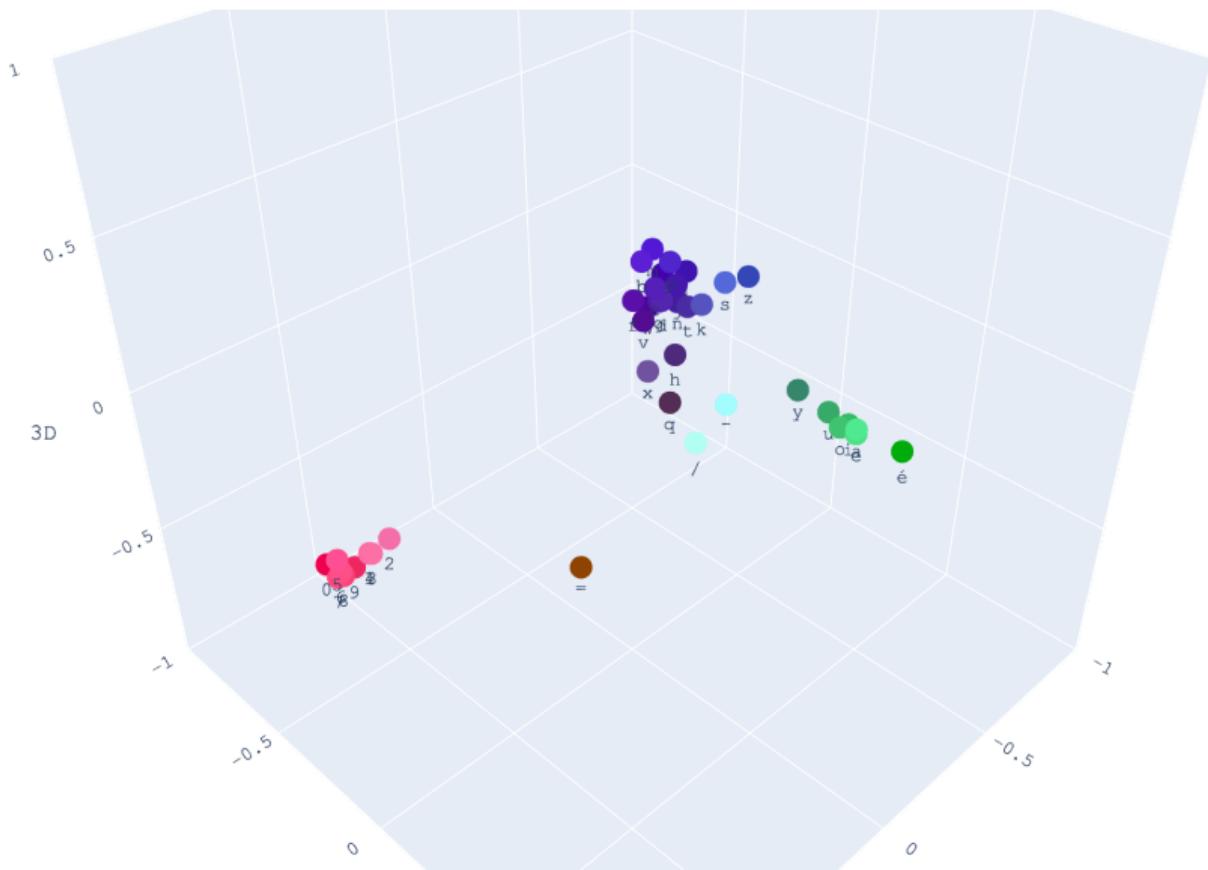
Truncate

$$\hat{U} \times \hat{\Sigma}$$

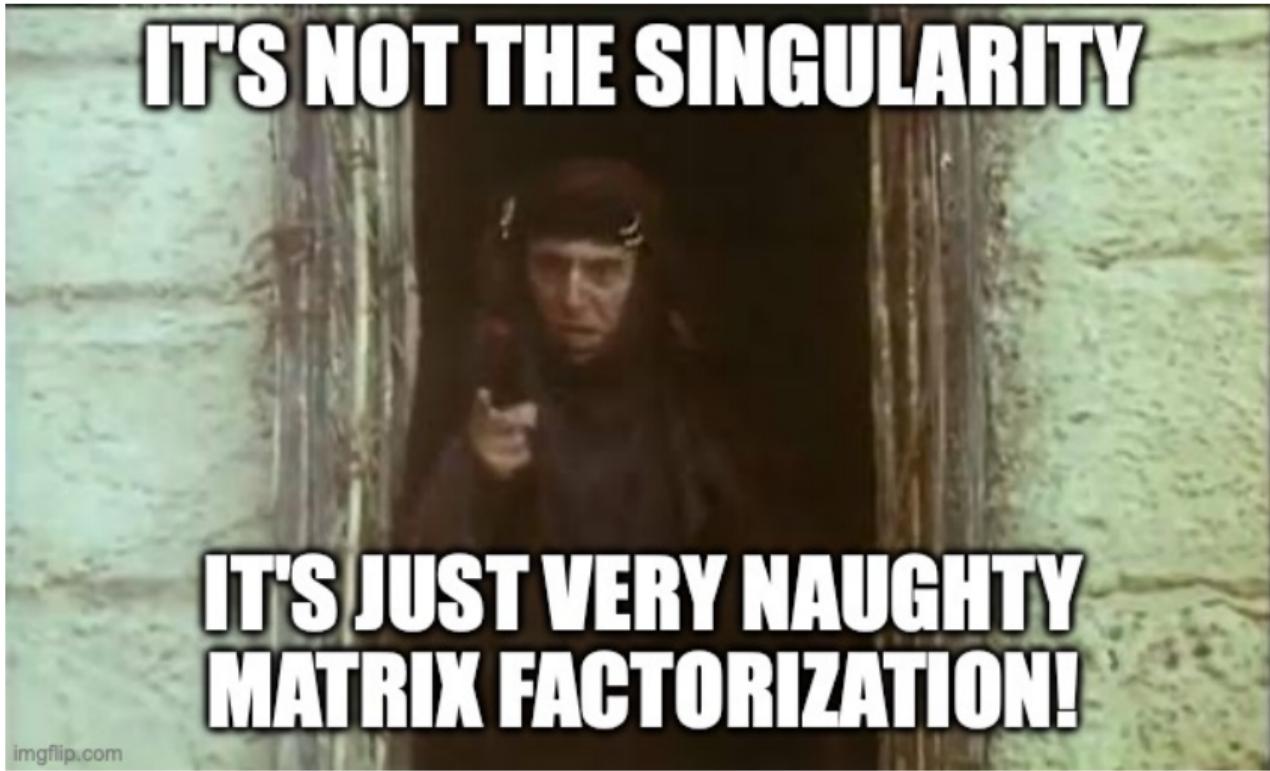


Plot

$$\hat{U} \times \hat{\Sigma}$$



What to conclude?

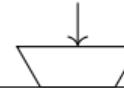


The Structure of Embeddings

Structure

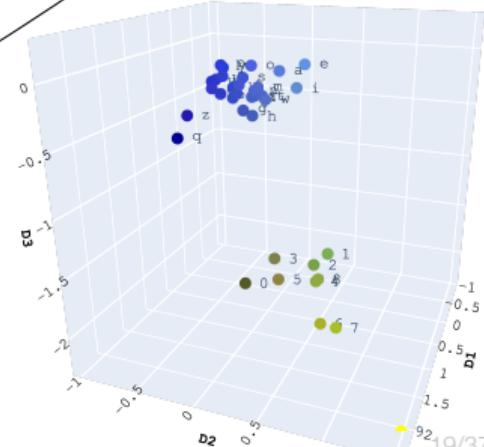
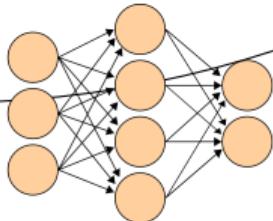
?

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



Embedding

Data

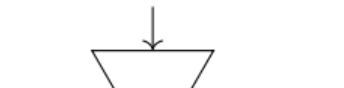


The Structure of Embeddings

Structure

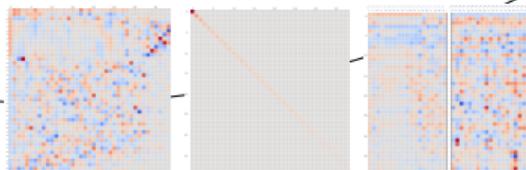


{-, /, 0, 1, 2, ..., 8, 9, =,
a, b, c, ..., w, x, y, z, é}

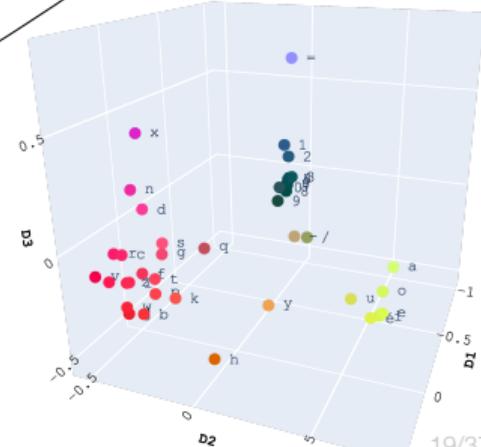


Embedding

Data



SVD



4 Why does this produce good word representations?

Good question. We don't really know.

The distributional hypothesis states that words in similar contexts have similar meanings. The objective above clearly tries to increase the quantity $v_w \cdot v_c$ for good word-context pairs, and decrease it for bad ones. Intuitively, this means that words that share many contexts will be similar to each other (note also that contexts sharing many words will also be similar to each other). This is, however, very hand-wavy.

Can we make this intuition more precise? We'd really like to see something more formal.

(Goldberg and Levy, 2014)

Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

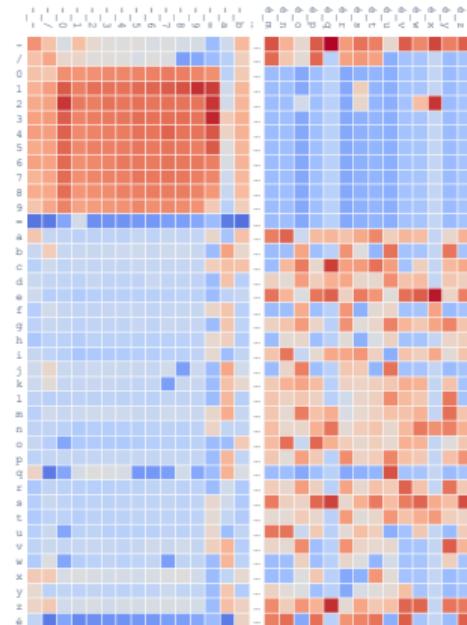
Embeddings as Functions Over Sets

$$\textcolor{red}{X} = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, \text{a}, \text{b}, \text{c}, \dots, \text{w}, \text{x}, \text{y}, \text{z}, \text{é}\}$$

$$\textcolor{blue}{Y} = X \times X = \{ (-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é}) \}$$

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

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



Embeddings as Functions Over Sets

$$\textcolor{red}{X} = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, \text{a}, \text{b}, \text{c}, \dots, \text{w}, \text{x}, \text{y}, \text{z}, \text{é}\}$$

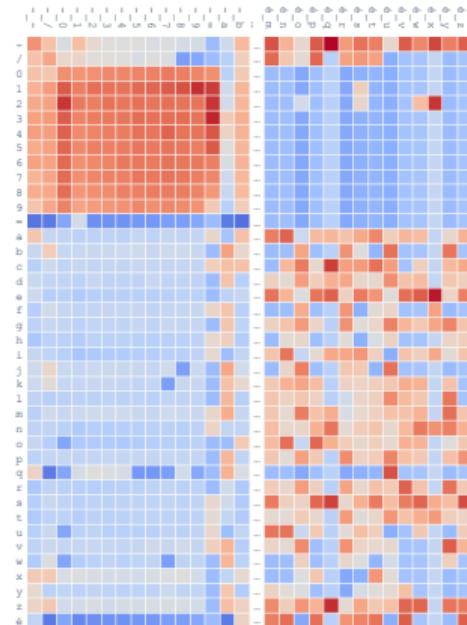
$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é})\}$$

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

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

$$M_x: \textcolor{red}{X} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

$$\textcolor{red}{x} \mapsto \textcolor{blue}{M}(x, -)$$



Embeddings as Functions Over Sets

$$\textcolor{red}{X} = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, \text{a}, \text{b}, \text{c}, \dots, \text{w}, \text{x}, \text{y}, \text{z}, \text{é}\}$$

$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é})\}$$

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

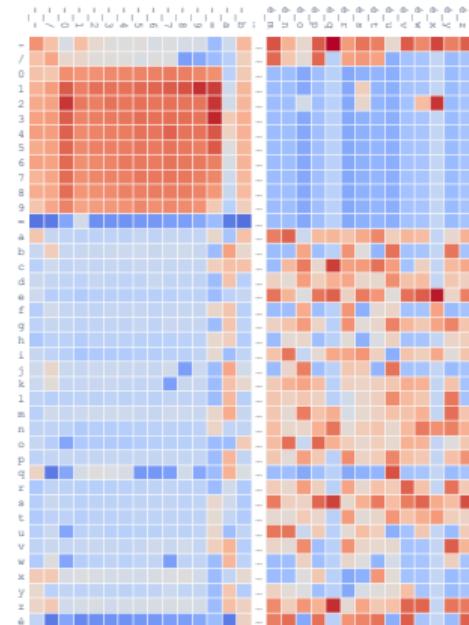
$$(\textcolor{red}{x}, \textcolor{blue}{y}) \mapsto \text{pmi}(\textcolor{red}{x}, \textcolor{blue}{y})$$

$$M_x: \textcolor{red}{X} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

$$\textcolor{red}{x} \mapsto M(x, -)$$

$$M_y: \textcolor{blue}{Y} \rightarrow \mathbb{R}^{\textcolor{red}{X}}$$

$$y \mapsto M(-, y)$$



Embeddings as Functions Over Sets

$$\textcolor{red}{X} = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, \text{a}, \text{b}, \text{c}, \dots, \text{w}, \text{x}, \text{y}, \text{z}, \text{é}\}$$

$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é})\}$$

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

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

$$\textcolor{red}{X} \xrightarrow{M_x} \mathbb{R}^{\textcolor{blue}{Y}}$$

$$M_x: \textcolor{red}{X} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

$$\textcolor{red}{x} \mapsto \textcolor{blue}{M}(x, -)$$

$$\mathbb{R}^{\textcolor{red}{X}} \xleftarrow{M_y} \textcolor{blue}{Y}$$

$$M_y: \textcolor{blue}{Y} \rightarrow \mathbb{R}^{\textcolor{red}{X}}$$

$$\textcolor{blue}{y} \mapsto \textcolor{red}{M}(-, y)$$

Embeddings as Functions Over Sets

$$\textcolor{red}{X} = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, \text{a}, \text{b}, \text{c}, \dots, \text{w}, \text{x}, \text{y}, \text{z}, \text{é}\}$$

$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é})\}$$

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

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

$$M_x: \textcolor{red}{X} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

$$\textcolor{red}{x} \mapsto \textcolor{blue}{M}(x, -)$$

$$M_y: \textcolor{blue}{Y} \rightarrow \mathbb{R}^{\textcolor{red}{X}}$$

$$\textcolor{blue}{y} \mapsto \textcolor{red}{M}(-, y)$$

$$\begin{array}{ccc} \textcolor{red}{X} & \xrightarrow{M_x} & \mathbb{R}^{\textcolor{blue}{Y}} \\ \downarrow & & \uparrow \\ \mathbb{R}^{\textcolor{red}{X}} & \xleftarrow{M_y} & \textcolor{blue}{Y} \end{array}$$

Embeddings as Functions Over Sets

$$\textcolor{red}{X} = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, \text{a}, \text{b}, \text{c}, \dots, \text{w}, \text{x}, \text{y}, \text{z}, \text{é}\}$$

$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é})\}$$

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

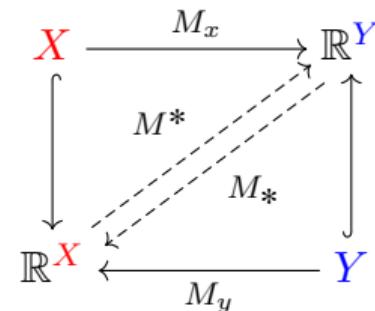
$$(\textcolor{red}{x}, \textcolor{blue}{y}) \mapsto \text{pmi}(\textcolor{red}{x}, \textcolor{blue}{y})$$

$$M_x: \textcolor{red}{X} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

$$\textcolor{red}{x} \mapsto \textcolor{blue}{M}(x, -)$$

$$M_y: \textcolor{blue}{Y} \rightarrow \mathbb{R}^{\textcolor{red}{X}}$$

$$y \mapsto \textcolor{red}{M}(-, y)$$

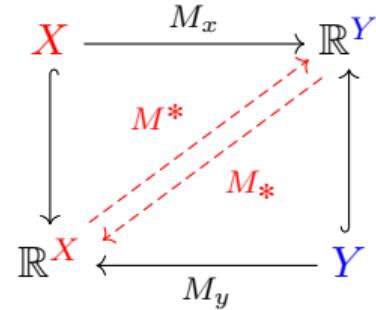


$$M^*: \mathbb{R}^{\textcolor{red}{X}} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

$$M_*: \mathbb{R}^{\textcolor{blue}{Y}} \rightarrow \mathbb{R}^{\textcolor{red}{X}}$$

Embeddings as Functions Over Sets

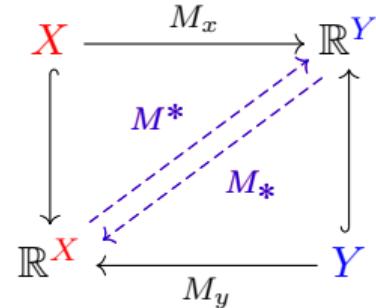
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$



Embeddings as Functions Over Sets

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

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



Embeddings as Functions Over Sets

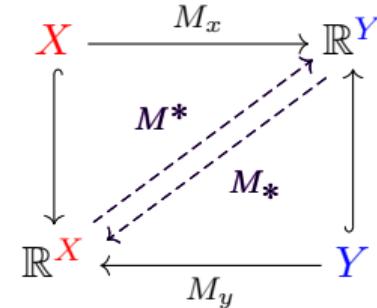
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

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

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

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

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



Embeddings as Functions Over Sets

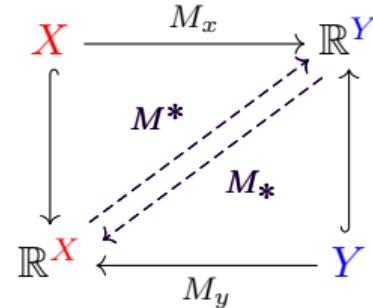
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

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

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

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

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



$$U := [\underline{u_1}, \dots, \underline{u_m}]$$

$$M = U \Sigma V^T \quad V := [\underline{v_1}, \dots, \underline{v_n}]$$

$$\Sigma := \begin{bmatrix} \sqrt{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sqrt{\lambda_r} \end{bmatrix}$$

Embeddings as Functions Over Sets

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

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

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

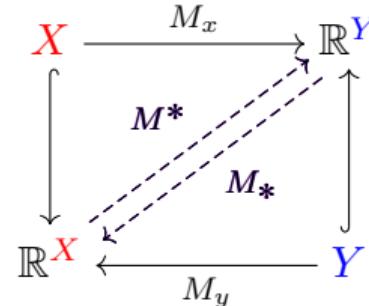
$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

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

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

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

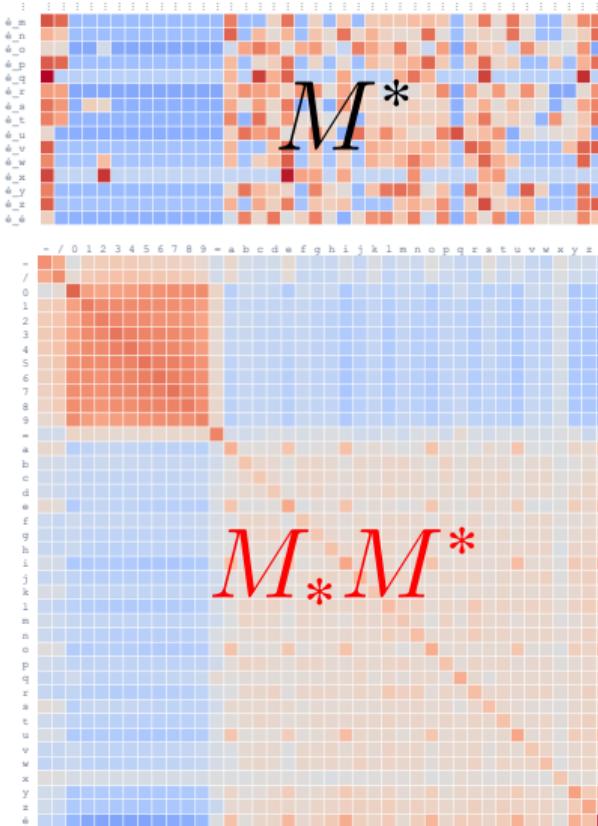
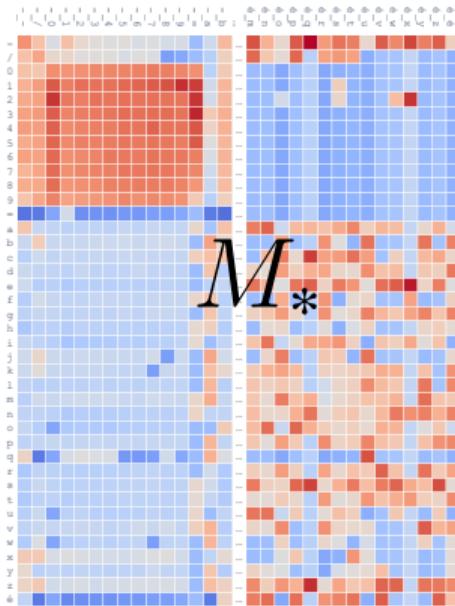
The u_i and v_i are (linear)
fixed points!



$$U := [u_1, \dots, u_m] \\ M = U \Sigma V^T \quad V := [v_1, \dots, v_n]$$

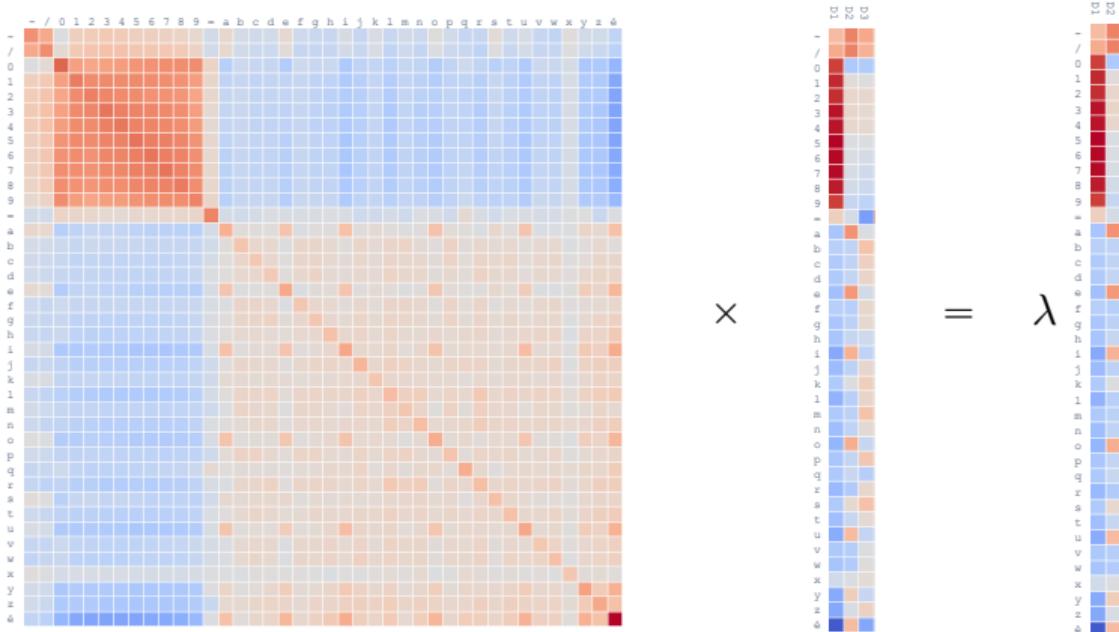
$$\Sigma := \begin{bmatrix} \sqrt{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sqrt{\lambda_r} \end{bmatrix}$$

$M_* M^*$ as a Covariance Matrix

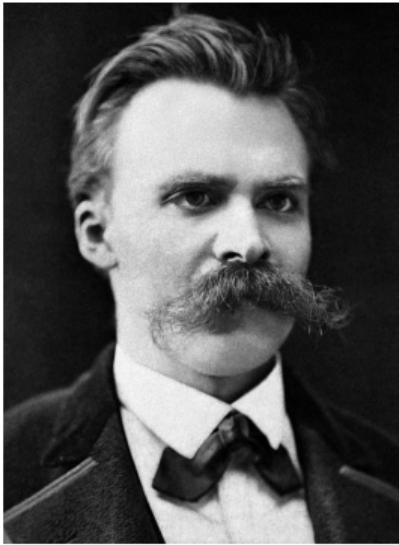


Eigenvectors as Fixed Points

$$M_* M^* u = \lambda u$$



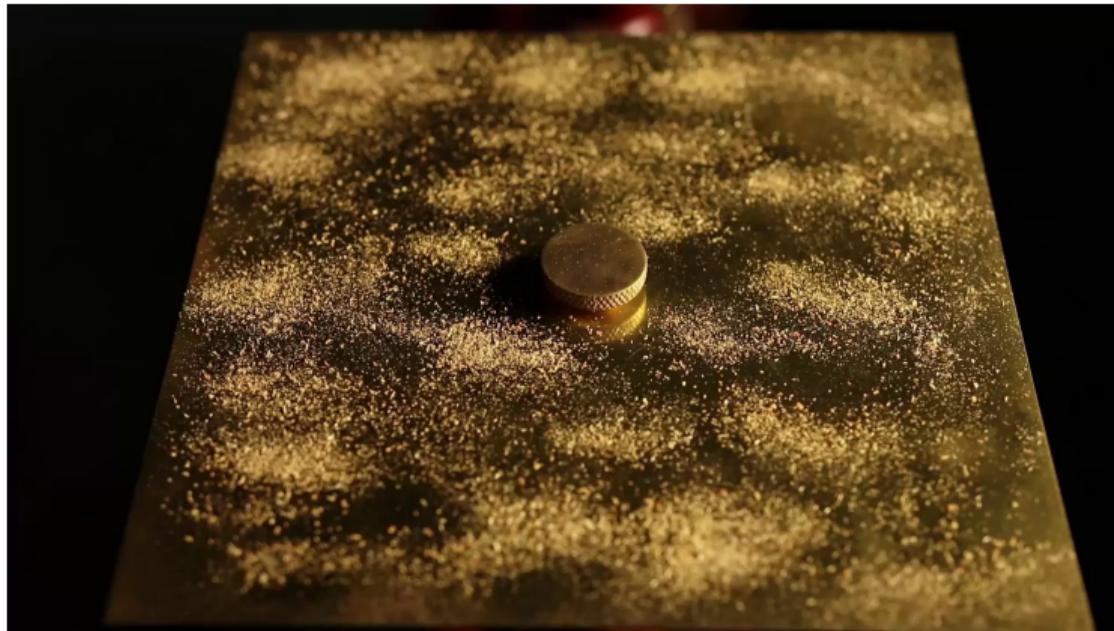
Chladni Figures



“One can conceive of a profoundly **deaf** human being who has never experienced sound or music; just as such a person will gaze in astonishment at the **Chladnian sound-figures in sand**, find their cause in the vibration of a string, and swear that **he must now know what men call sound** — this is precisely what happens to all of us with **language**.[”]

(Nietzsche, 1873)

Chladni Figures

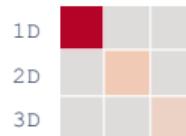


Structural Features

Eigenvectors of $M_* M^*$:



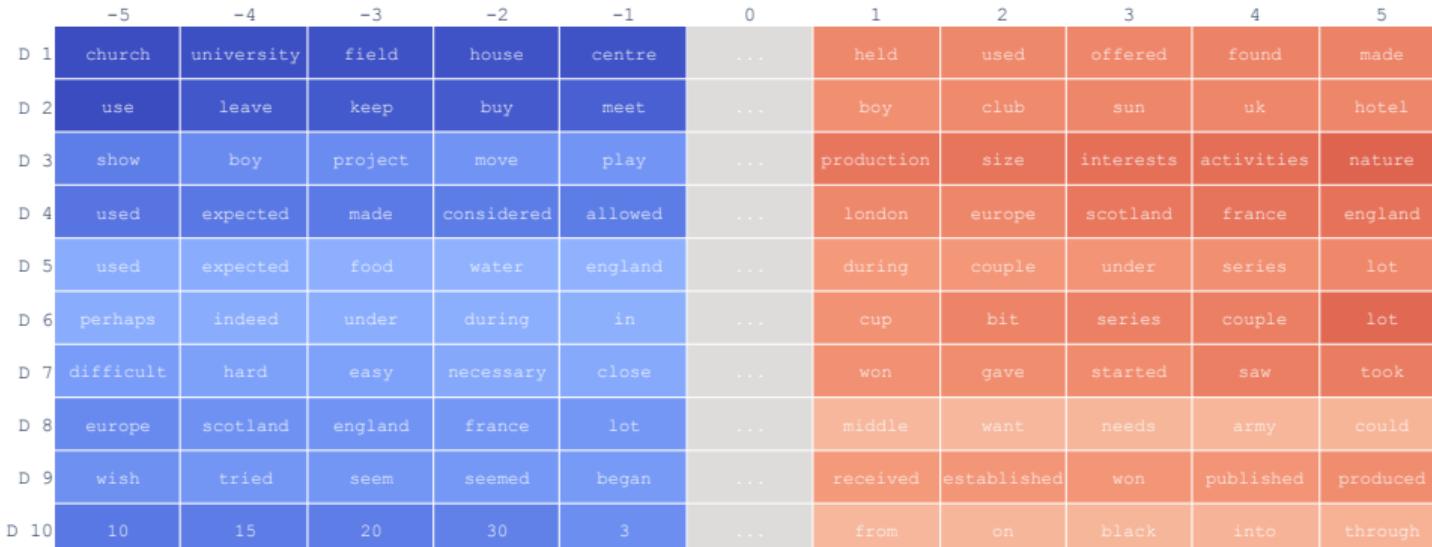
Eigenvalues of $M_* M^*$ and $M^* M_*$:



Eigenvectors of $M^* M_*$:



Words



Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

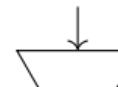
Conclusion

The Structure of Embeddings

Structure

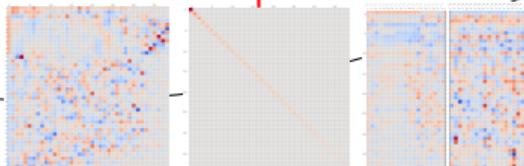


$\{-, /, 0, 1, 2, \dots, 8, 9, =,$
 $a, b, c, \dots, w, x, y, z, \acute{e}\}$

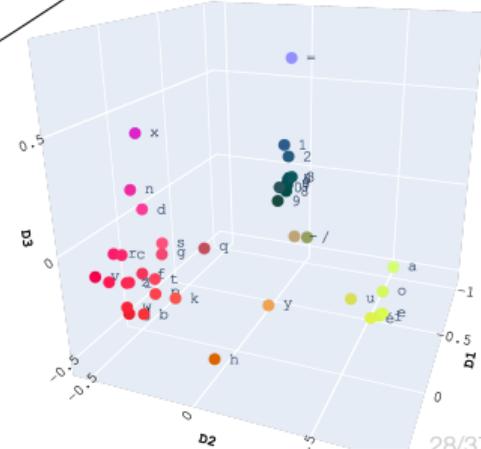


Embedding

Data



SVD

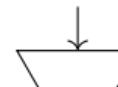


The Structure of Embeddings

Structure

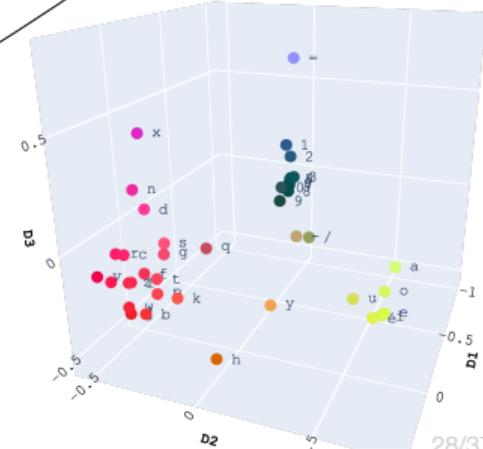
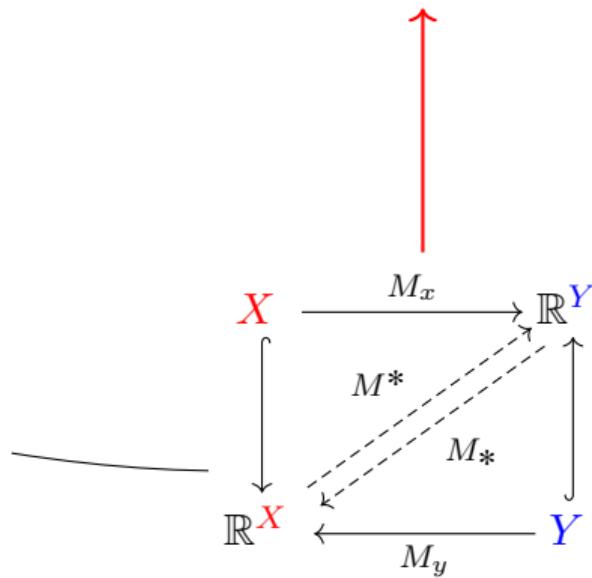


$\{-, /, 0, 1, 2, \dots, 8, 9, =,$
 $a, b, c, \dots, w, x, y, z, \acute{e}\}$



Embedding

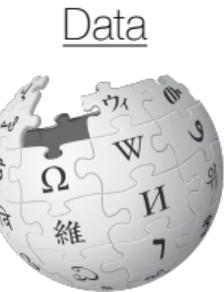
Data



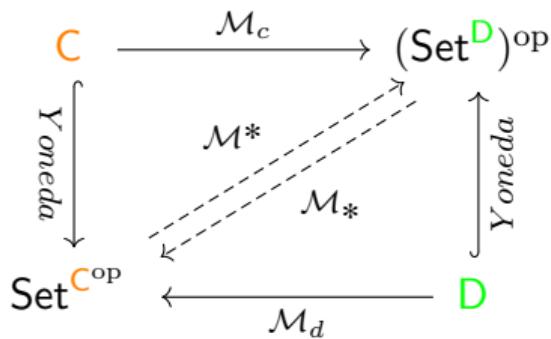
The Structure of Embeddings

Structure

?

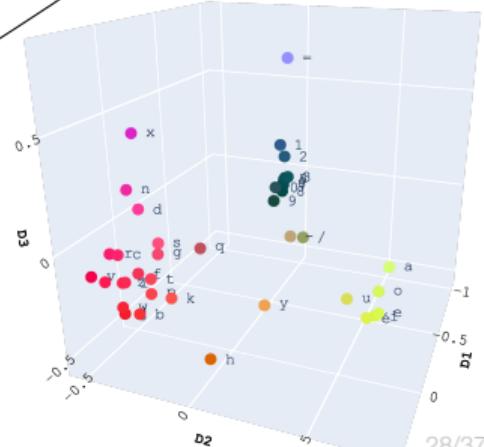
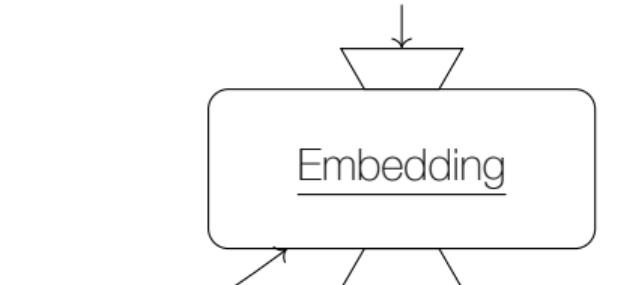


Data



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

Embedding

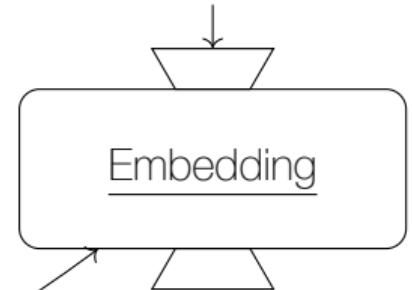


The Structure of Embeddings

Structure

?

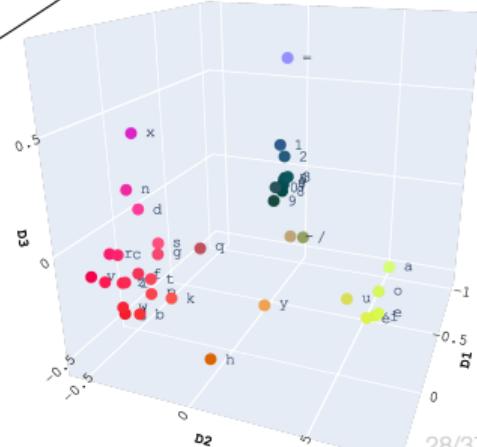
$\{-, /, 0, 1, 2, \dots, 8, 9, =,$
 $a, b, c, \dots, w, x, y, z, é\}$



Data



$$\textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \text{Set}$$



Structure

?

$$\begin{array}{ccc} \textcolor{teal}{term}_i & \textcolor{teal}{context}_i & \text{measure} \\ \searrow & \downarrow & \swarrow \\ \textcolor{orange}{C}^{\text{op}} & \times \textcolor{green}{D} & \rightarrow \mathbf{Set} \end{array}$$

Structure

?

$$\begin{array}{ccc} \textcolor{teal}{term}_i & \textcolor{teal}{context}_i & \text{measure} \\ \searrow & \downarrow & \swarrow \\ \textcolor{orange}{C}^{\text{op}} & \times \textcolor{green}{D} & \rightarrow \textcolor{red}{Set} \end{array}$$

Structure

?

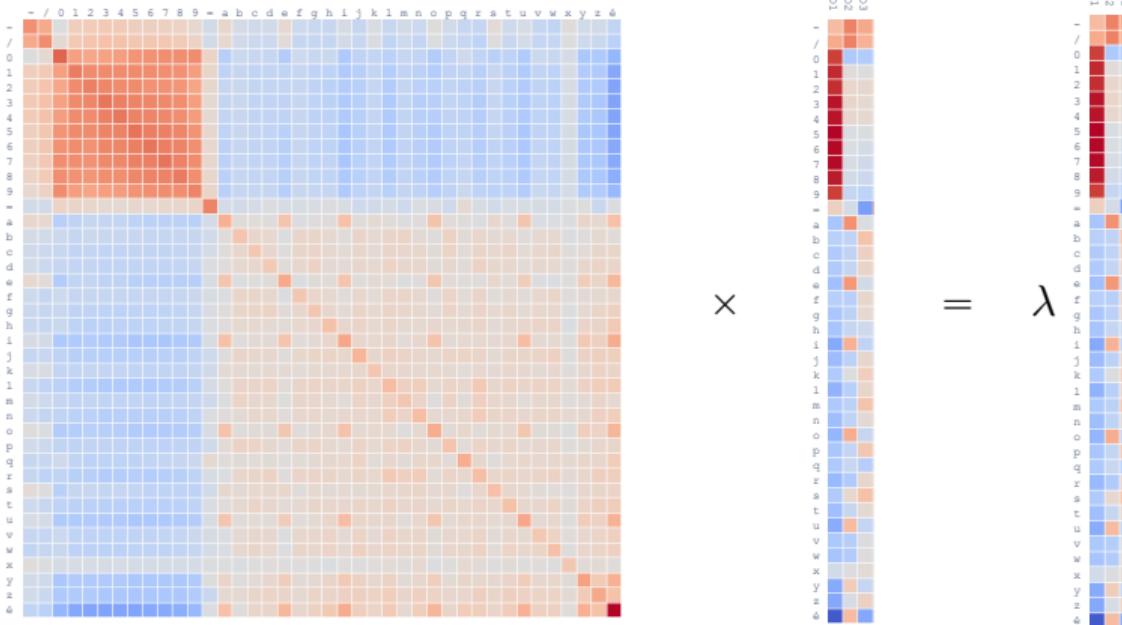
$$\begin{array}{ccc} \textcolor{teal}{term}_i & \textcolor{teal}{context}_i & \text{measure} \\ \downarrow & \downarrow & \swarrow \\ \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \textcolor{red}{2} \end{array}$$

Structure

$$\begin{array}{c} \text{C}^{\text{op}} \times \text{D} \rightarrow 2 \\ \Downarrow \\ \mathcal{M}^*: 2^{\text{C}^{\text{op}}} \rightleftarrows (2^{\text{D}})^{\text{op}}: \mathcal{M}_* \end{array}$$

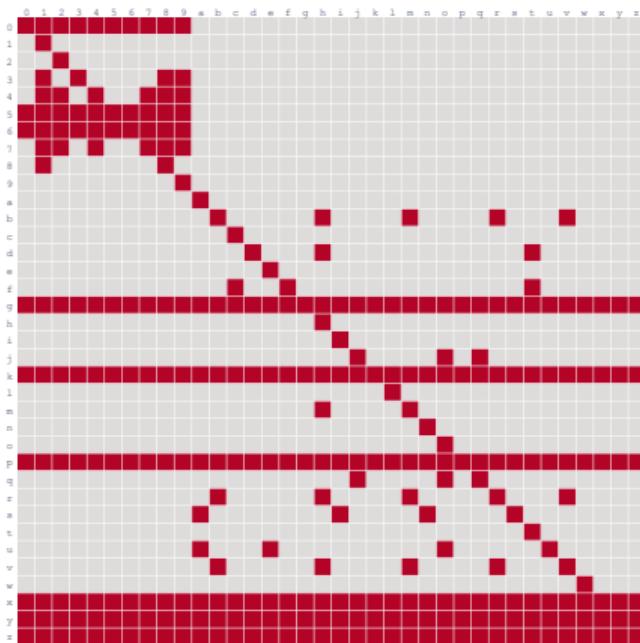
Binary Fixed Points

$$M_* M^* u = \lambda u$$



Binary Fixed Points

$$\mathcal{M}_*\mathcal{M}^*f = f$$



★

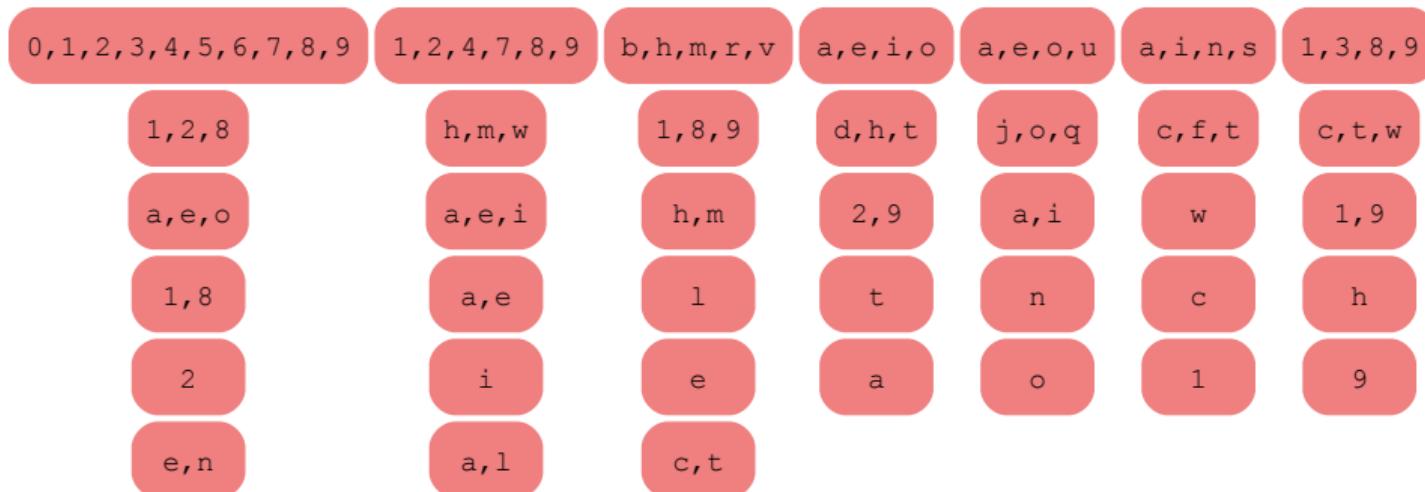


?

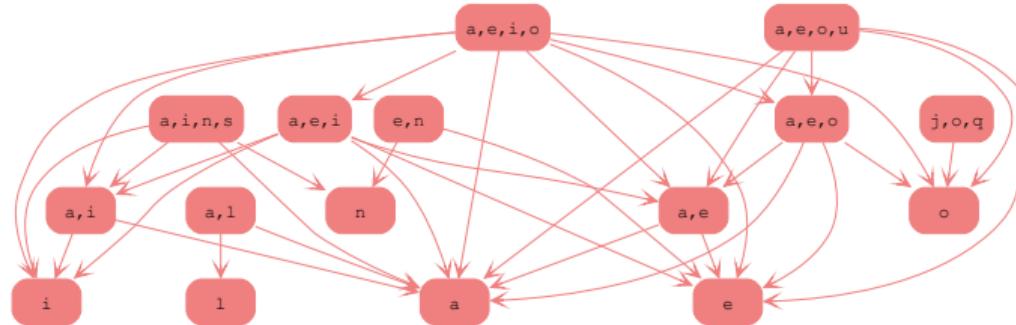
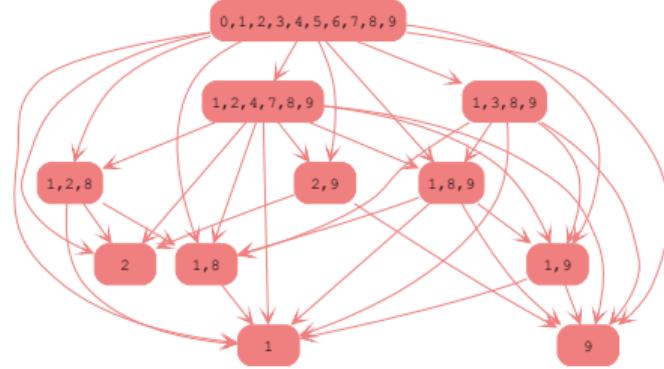
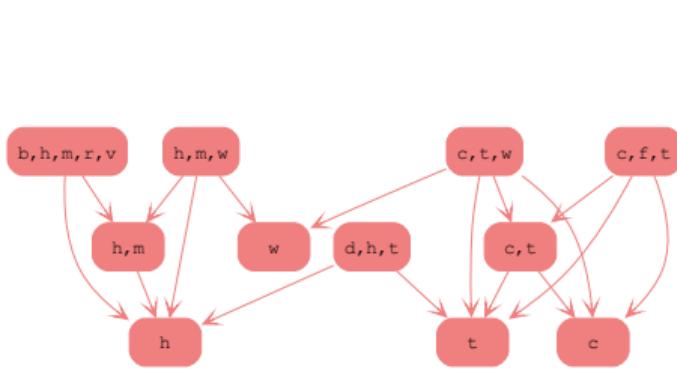


“Eigensets”

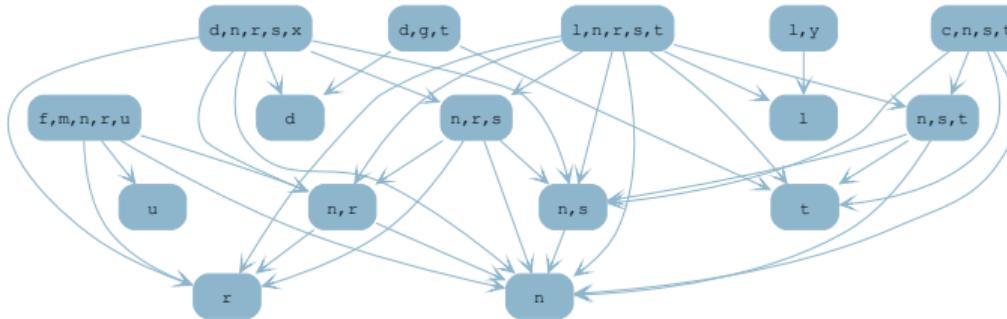
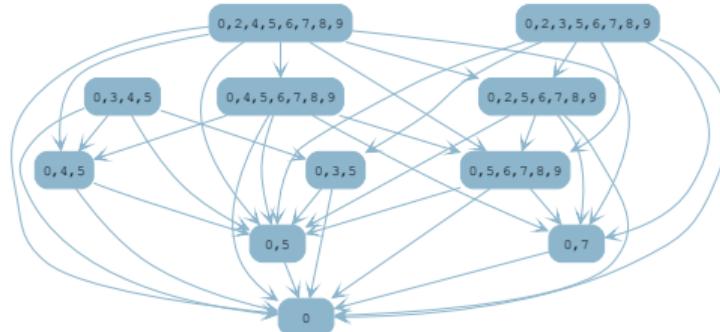
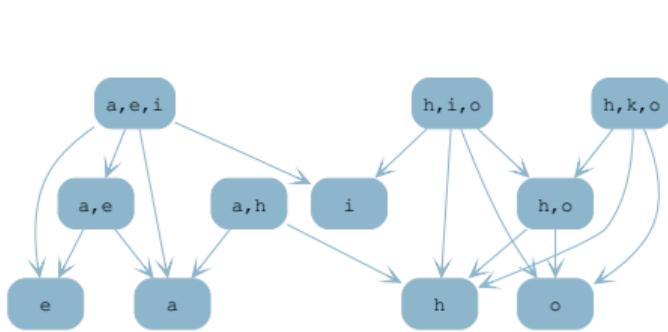
$$\mathcal{M}_*\mathcal{M}^*f = f$$



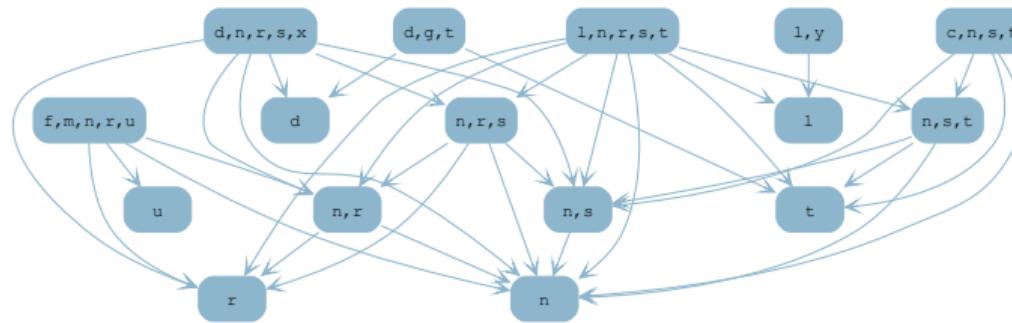
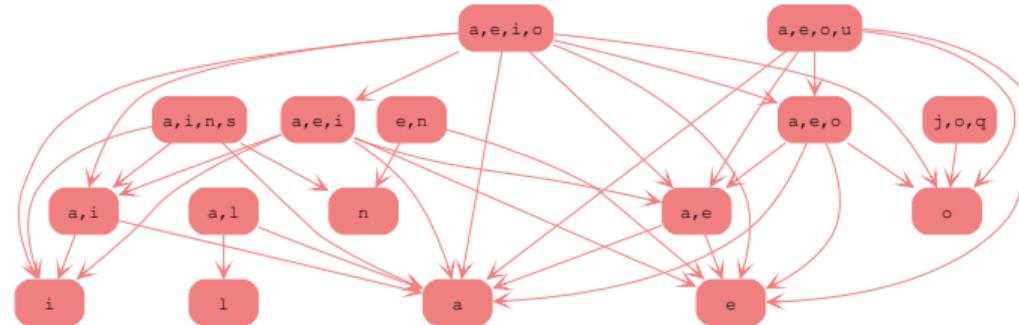
Partial Order Structure

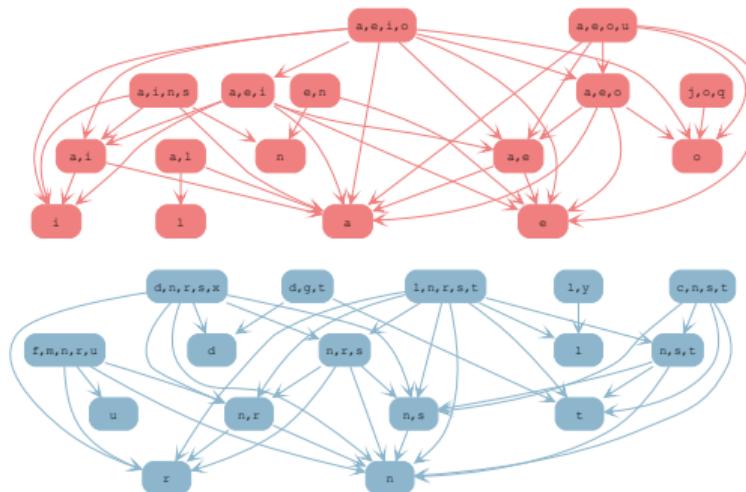


Dual Partial Order



Paring of Partial Ordered Fixed Points



Structure

$$\begin{array}{c} \text{C}^{\text{op}} \times \text{D} \rightarrow 2 \\ \Downarrow \\ \mathcal{M}^*: 2^{\text{C}^{\text{op}}} \rightleftarrows (2^{\text{D}})^{\text{op}}: \mathcal{M}_* \end{array}$$

Structure

?

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

Structure

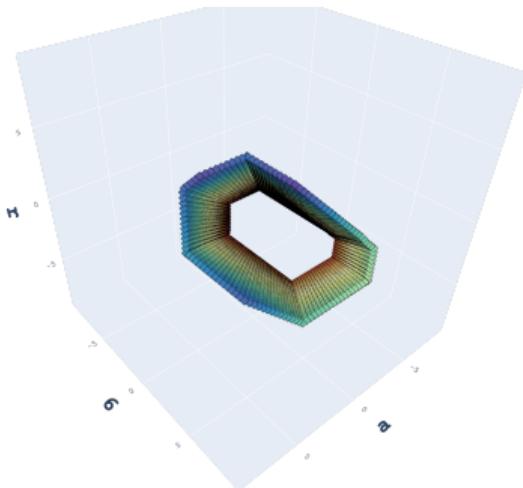
?



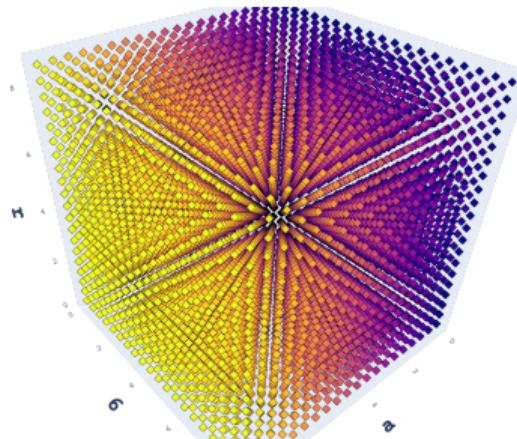
$$\begin{array}{c} \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}} \\ \Downarrow \\ \mathcal{M}^*: \bar{\mathbb{R}}^{\textcolor{orange}{C}^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^{\textcolor{green}{D}})^{\text{op}}: \mathcal{M}_* \end{array}$$

Enriching over $\bar{\mathbb{R}}$

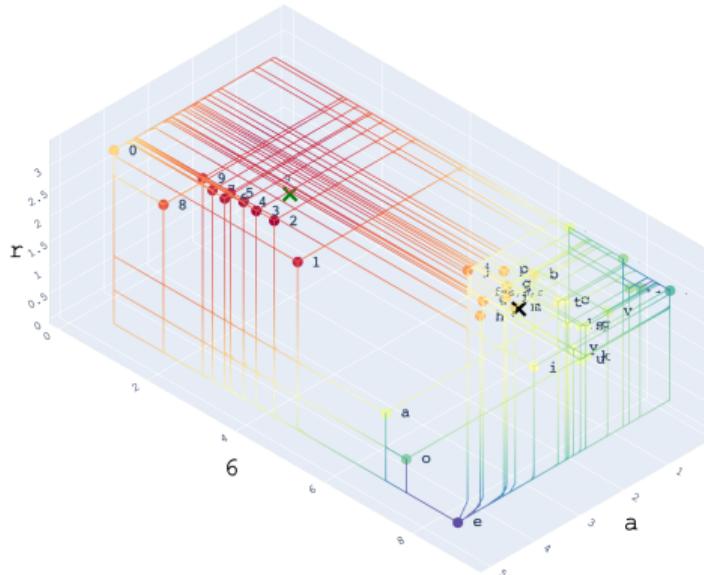
Structure



$$\leftarrow \mathcal{M}_* \mathcal{M}^*$$



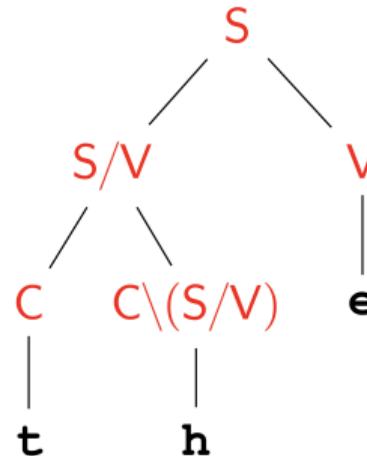
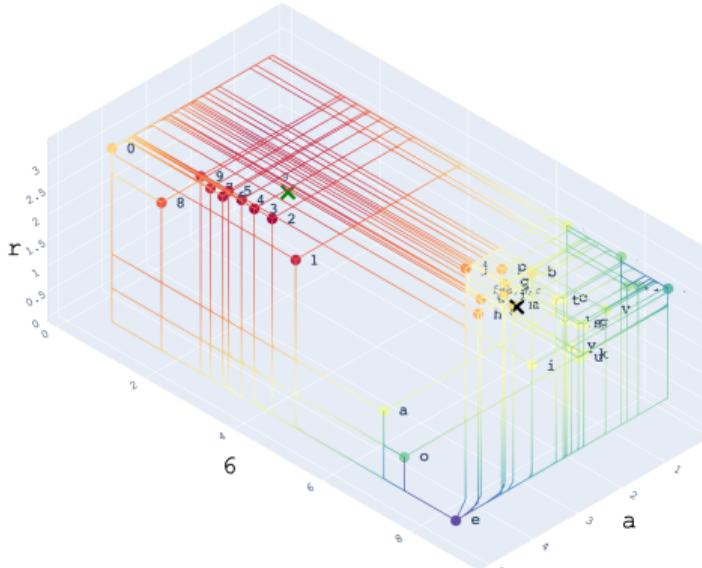
$$\begin{array}{c} \text{C}^{\text{op}} \times \text{D} \rightarrow \bar{\mathbb{R}} \\ \Downarrow \\ \mathcal{M}^*: \bar{\mathbb{R}}^{\text{C}^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^{\text{D}})^{\text{op}}: \mathcal{M}_* \end{array}$$

Structure

$$\begin{array}{c}
 \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}} \\
 \Downarrow \\
 \mathcal{M}^*: \bar{\mathbb{R}}^{\textcolor{orange}{C}^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^{\textcolor{green}{D}})^{\text{op}}: \mathcal{M}_*
 \end{array}$$

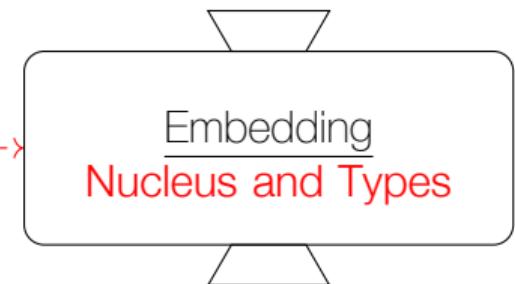
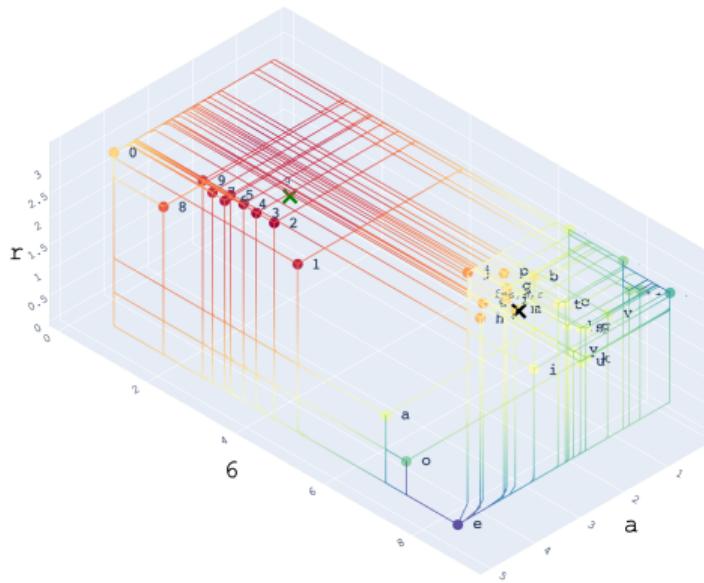
Enriching over $\bar{\mathbb{R}}$

Structure



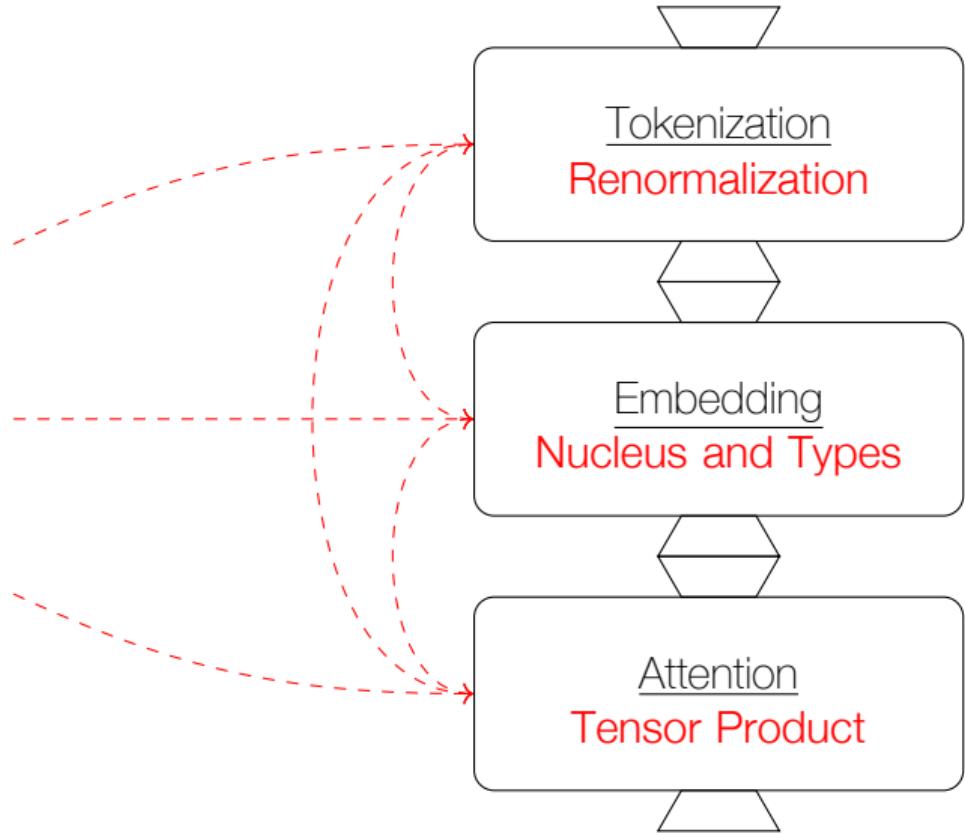
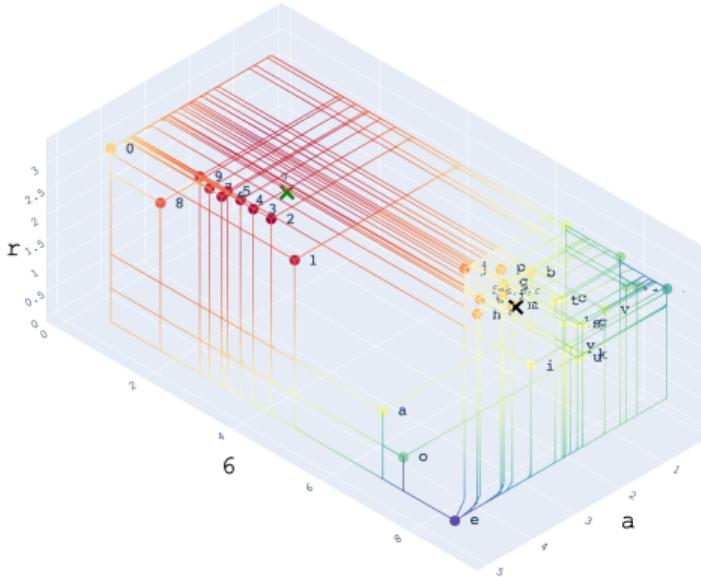
$$\begin{array}{c}
 \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}} \\
 \Downarrow \\
 \mathcal{M}^*: \bar{\mathbb{R}}^{\textcolor{orange}{C}^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^{\textcolor{green}{D}})^{\text{op}}: \mathcal{M}_*
 \end{array}$$

Structure



Formal Explainability

Structure



Outline

Intro: Critique and Formalism

Epistemological Critique: LLMs as Formal Objects

Theoretical Critique: Formal Explainability

The Algebra Behind the Embeddings

The Structure Behind the Algebra

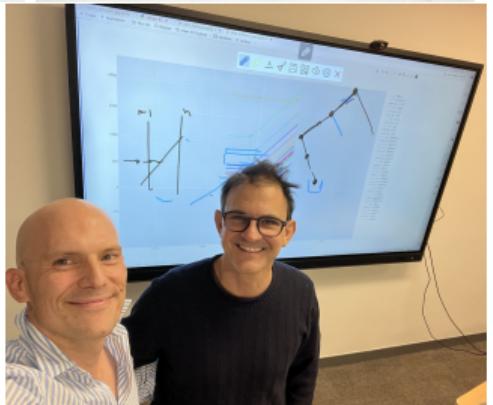
The Categories Behind the Structure

Conclusion

Conclusion: For a Critical Formalism

- ◊ It is urgent to address the **epistemological** dimension of the critical project in its own terms
- ◊ This requires to develop a **critical approach within formal sciences** where formalization is not assumed to lead to **naturalization**
 - ◊ The new role of **data** within formal sciences is crucial in this sense
- ◊ A **critical formalism** will be incomplete if it remains disconnected from the **political**, and even the **artistic** dimension of the critical program
 - ◊ Pure form of data articulated with a **politics of the corpus**
 - ◊ We need a **new alliance** between the **formal sciences**, the **human sciences**, and the **arts**.

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*. <https://api.semanticscholar.org/CorpusID:263613625>

Références 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>
- Bender, E. M., & Koller, A. (2020). Climbing towards NLU: On meaning, form, and understanding in the age of data. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5185–5198. <https://doi.org/10.18653/v1/2020.acl-main.463>
- Bourdieu, P. (1979). *La distinction: Critique sociale du jugement*. Éditions de Minuit.
- Bourdieu, P. (1994). *Raisons pratiques: Sur la théorie de l'action*. Éditions du Seuil.
- 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*. <https://api.semanticscholar.org/CorpusID:263613625>
- Church, A. (1936). An unsolvable problem of elementary number theory. *American Journal of Mathematics*, 58(2), 345–363.
- Foucault, M. (1966). *Les mots et les choses : Une archéologie des sciences humaines*. Gallimard.
- 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>
- Girard, J.-Y. (2006). *Le point aveugle: Cours de logique. vers la perfection*. Editions Hermann.
- Gödel, K. (1934). On undecidable propositions of formal mathematical systems. In *Collected works* (pp. 346–371). Clarendon Press Oxford University Press.

Références II

- Goldberg, Y., & Levy, O. (2014). Word2vec explained: Deriving mikolov et al.'s negative-sampling word-embedding method. *CoRR*, abs/1402.3722.
- Harris, Z. (1960). *Structural linguistics*. University of Chicago Press.
- Hjelmslev, L. (1935). *La catégorie des cas*. Wilhelm Fink Verlag.
- Hjelmslev, L. (1971). La structure fondamentale du langage. In *Prolégomènes à une théorie du langage* [Prolégomènes à une theorie du langage] (pp. 177–231). Éditions de Minuit.
- Hjelmslev, L. (1975). *Résumé of a Theory of Language*. Nordisk Sprog-og Kulturforlag.
- Jakobson, R., Fant, G. M., & Halle, M. (1952). *Preliminaries to speech analysis: The distinctive features and their correlates*. MIT Press.
- 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/>
- Latour, B., Jensen, P., Venturini, T., Grauwin, S., & Boullier, D. (2012). 'The whole is always smaller than its parts' - a digital test of Gabriel Tardes' monads. *The British Journal of Sociology*, 63(4), 590–615.
- Lévi-Strauss, C. (1949). *Les structures élémentaires de la parenté*. Presses Universitaires de France.
- Lévi-Strauss, C. (1962). *La pensée sauvage*. Plon.
- 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.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *CoRR*, abs/1310.4546.

Références III

- Nietzsche, F. (1873). On truth and lying in a non-moral sense [Originally unpublished; written in 1873.]. (R. Speirs, Trans.). In R. Geuss & R. Speirs (Eds.), *The birth of tragedy and other writings* (pp. 141–153). Cambridge University Press.
- 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.
- Spang-Hanssen, H. (1959). *Probability and structural classification in language description*. Rosenkilde; Bagger.
- Turing, A. (1937). On Computable Numbers, with an Application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, s2-42(1), 230–265. <https://doi.org/10.1112/plms/s2-42.1.230>
- Underwood, T. (2023, October 15). *The empirical triumph of theory* [Accessed: 2023-10-15].
<https://critinq.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

Chat Token Vector
Università Ca' Foscari
Venice, Italy

Toward a Critical Formalism

Philosophical and Theoretical Effects of a Mathematical Critique of LLMs

Juan Luis Gastaldi

www.giannigastaldi.com

ETH zürich

June 12, 2025