

*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](http://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)

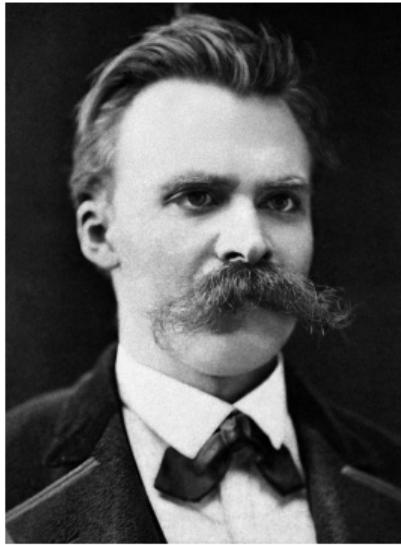
# Where Art Thou, Critique?

- ◊ 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.”

# 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?]

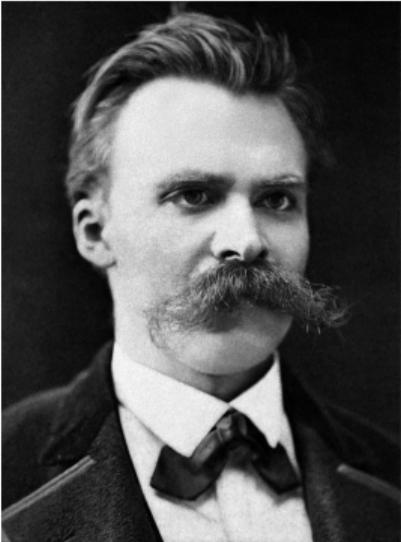
Any regularity in language/knowledge is not natural but cultural/social/political  
(Political)

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

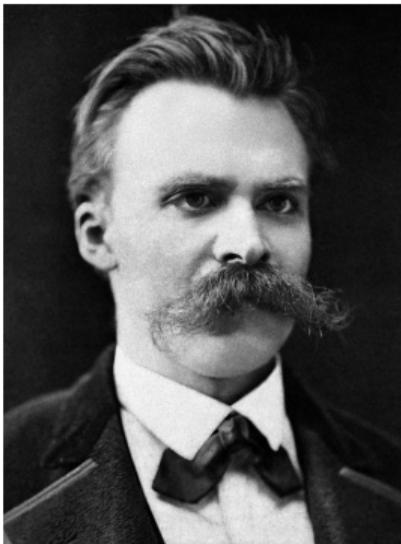
- ◊ At the source of this situation is the new foundational role played by **formal sciences** in the 20th century
  - ◊ For a **theory of language**: Carnap, Gödel, Turing, Shannon, Harris, Chomsky...

- ◊ 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

# Critique vs Formalism

- 
- ◊ “their eyes merely glide across the surface of things and see ‘forms’; nowhere does their perception lead into truth; instead it is content to receive stimuli and, as it were, to play with its fingers on the back of things.”

# Critique vs Formalism

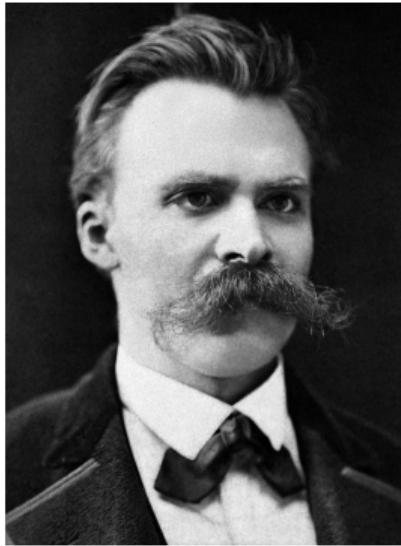


- ◊ “their eyes merely glide across the surface of things and see ‘forms’; nowhere does their perception lead into truth; instead it is content to receive stimuli and, as it were, to play with its fingers on the back of things.”
- ◊ “...the great edifice of concepts exhibits the rigid regularity of a Roman columbarium, while logic breathes out that air of severity and coolness which is peculiar to mathematics.”
- ◊ “...if we are forced to comprehend all things under these forms alone, then it is no longer wonderful that what we comprehend in all these things is actually nothing other than these very forms; for all of them must exhibit the laws of number, and number is precisely that which is most astonishing about things.”
- ◊ “...the only things we really know about them [laws of nature] are things which we bring to bear on them: time and space, in other words, relations of succession and number.”

- ◊ 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**

# For a Critical Formalism



“...all peoples have just such a **mathematically divided firmament of concepts** above them [...] Here, one can certainly **admire humanity** as a mighty **architectural genius** who succeeds in erecting the infinitely complicated cathedral of concepts **on moving foundations**, or even, one might say, **on flowing water**; admittedly, in order to rest on such foundations, it has to be like a thing constructed from cobwebs, **so delicate** that it can be carried off on the waves and yet **so firm** as not to be blown apart by the wind.”

(Nietzsche, 1873)

- ◊ 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



?

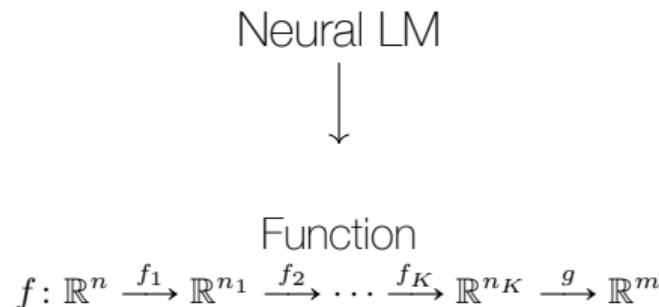
# Neural LMs as Computable Functions

Neural LM

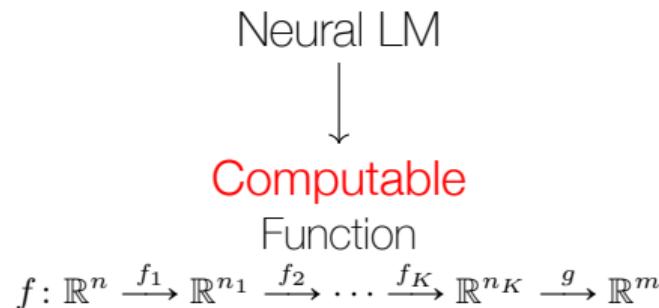


$f$  !

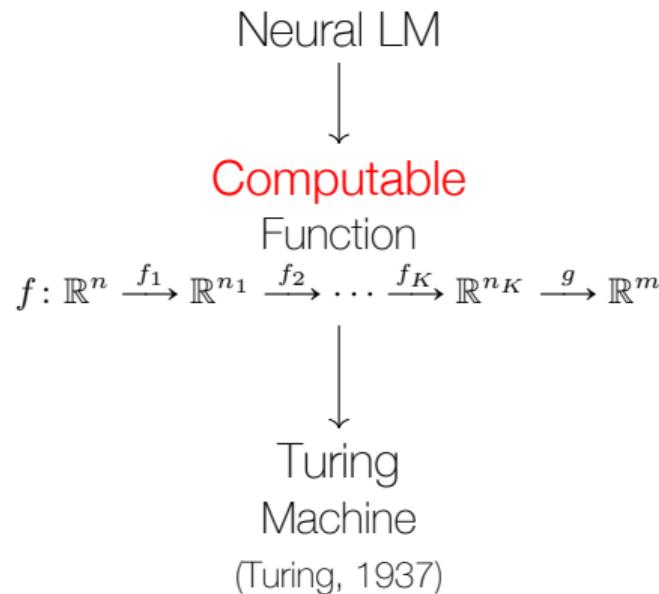
# Neural LMs as Computable Functions



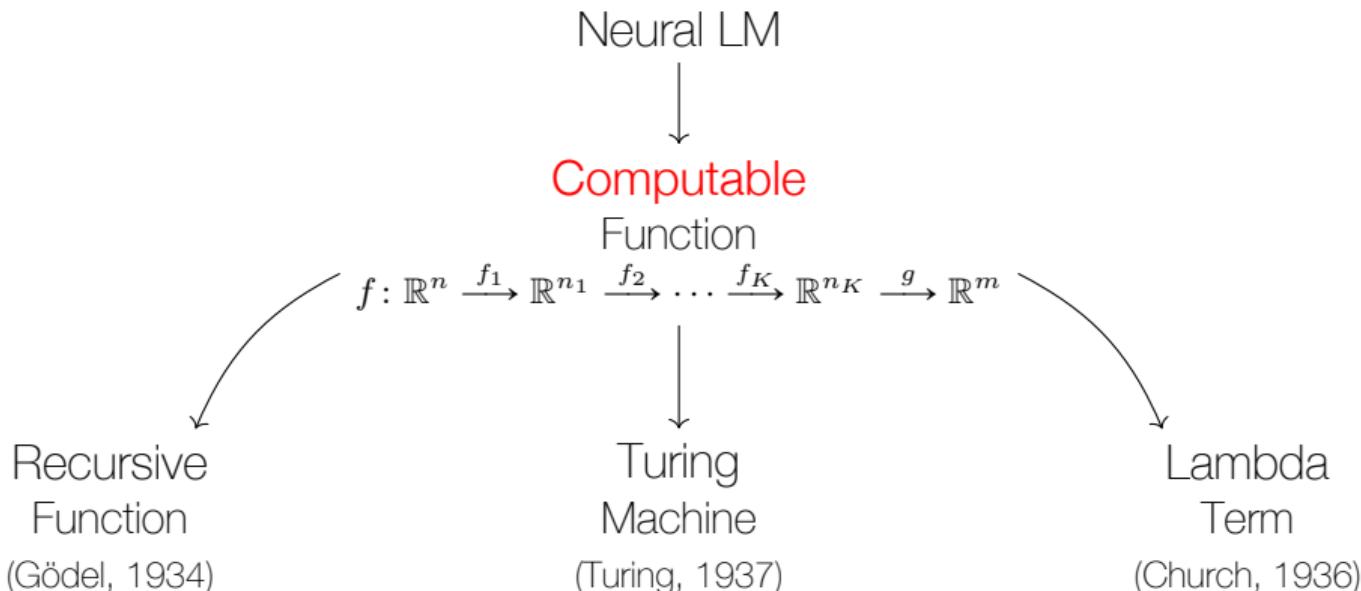
# Neural LMs as Computable Functions



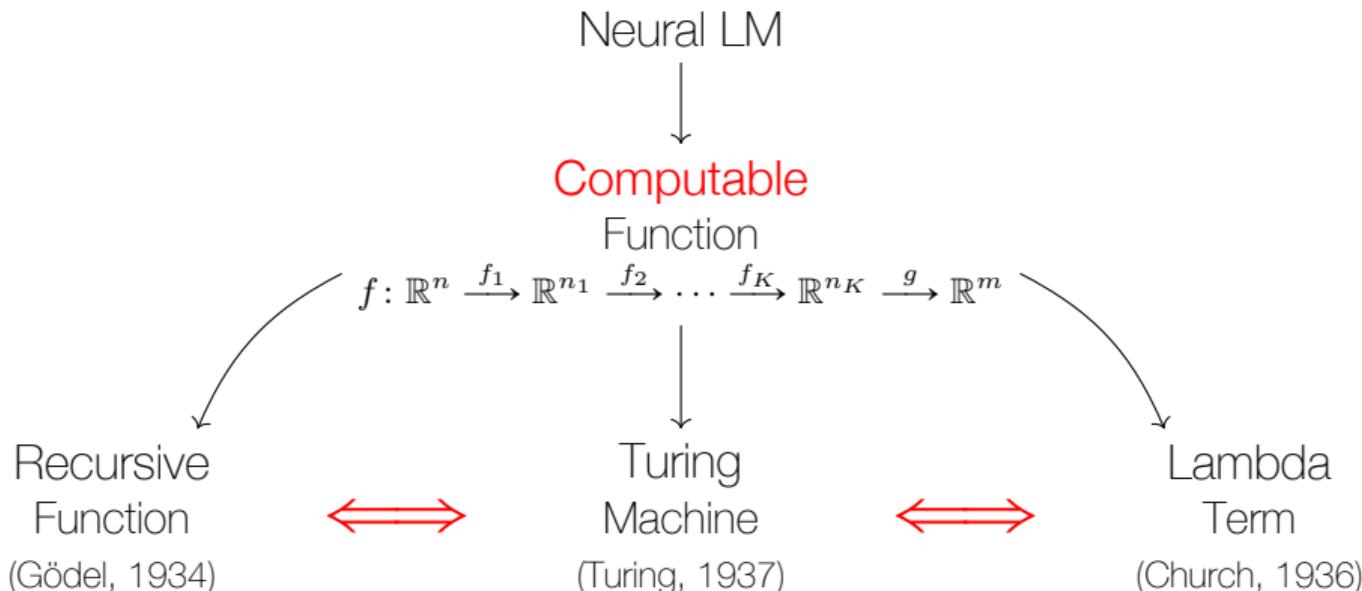
# Neural LMs as Computable Functions



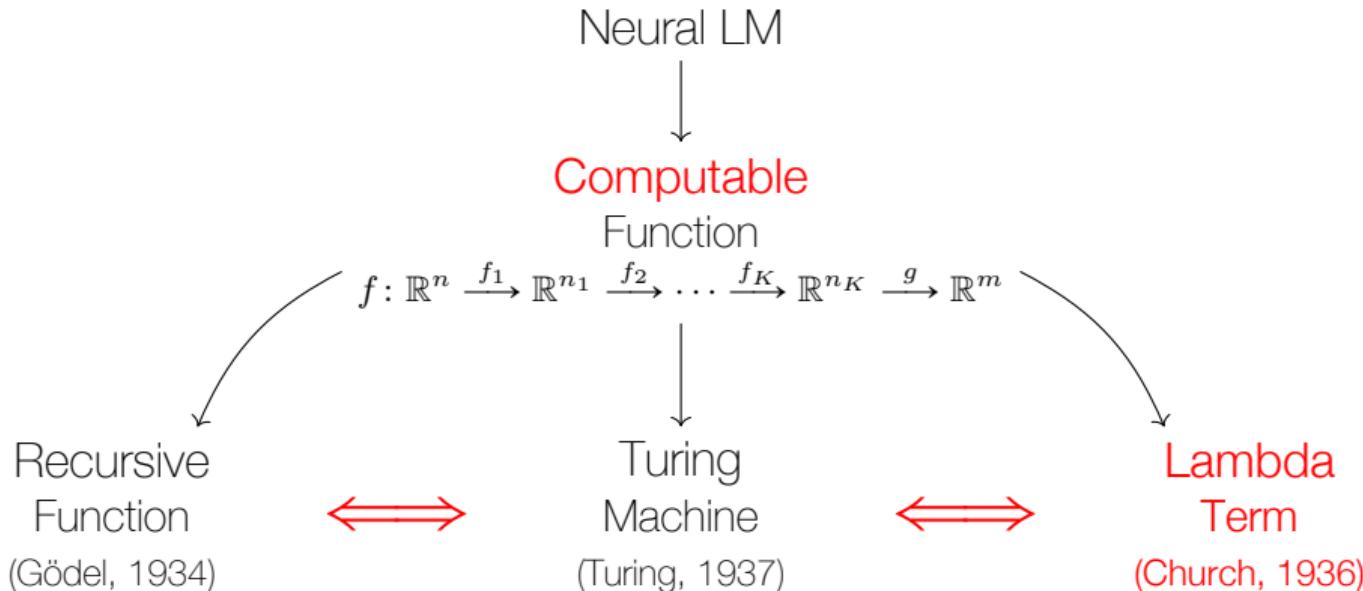
# Neural LMs as Computable Functions



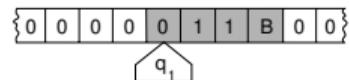
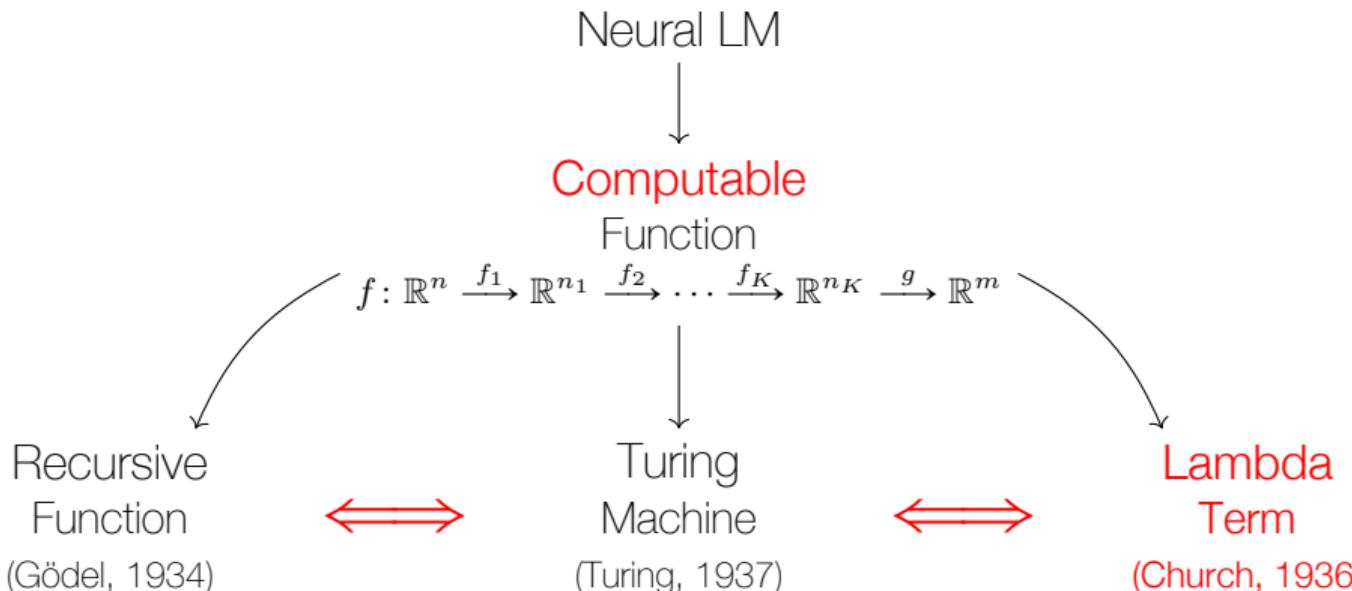
# Neural LMs as Computable Functions



# 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

# $\lambda$ -abstraction and $\beta$ -reduction in $\lambda$ -calculus

$yxz$

# $\lambda$ -abstraction and $\beta$ -reduction in $\lambda$ -calculus

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

# $\lambda$ -abstraction and $\beta$ -reduction in $\lambda$ -calculus

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

# $\lambda$ -abstraction and $\beta$ -reduction in $\lambda$ -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:  $\lambda f. \lambda x. f(\underbrace{\dots (f x)}_{n \text{ times}} \dots)$

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

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

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

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

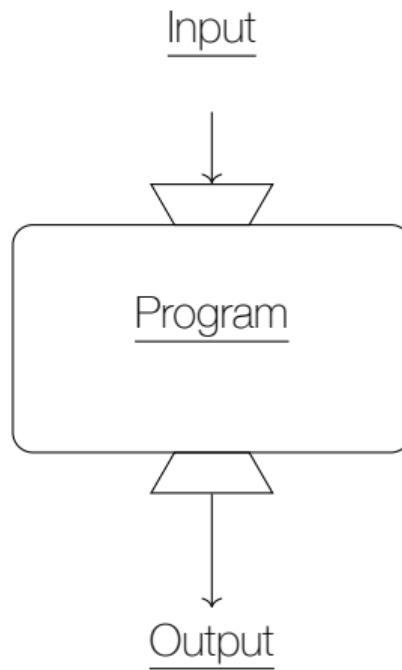
$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f(n f x)$$

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:	$\lambda f. \lambda x. f(f(f(f(fx))))$	$\lambda f. \lambda x. (\lambda y. f(fy)) ((\lambda h. \lambda z. h(h(hz))) f x)$
...		$\lambda f. \lambda x. (\lambda y. f(fy)) ((\lambda z. f(f(fz))) x)$
$n:$	$\lambda f. \lambda x. \underbrace{f(\dots(f}_{n \text{ times}} x) \dots)$	$\lambda f. \lambda x. (\lambda y. f(fy)) (f(f(fx)))$
		$\lambda f. \lambda x. f(f(f(f(fx)))))$

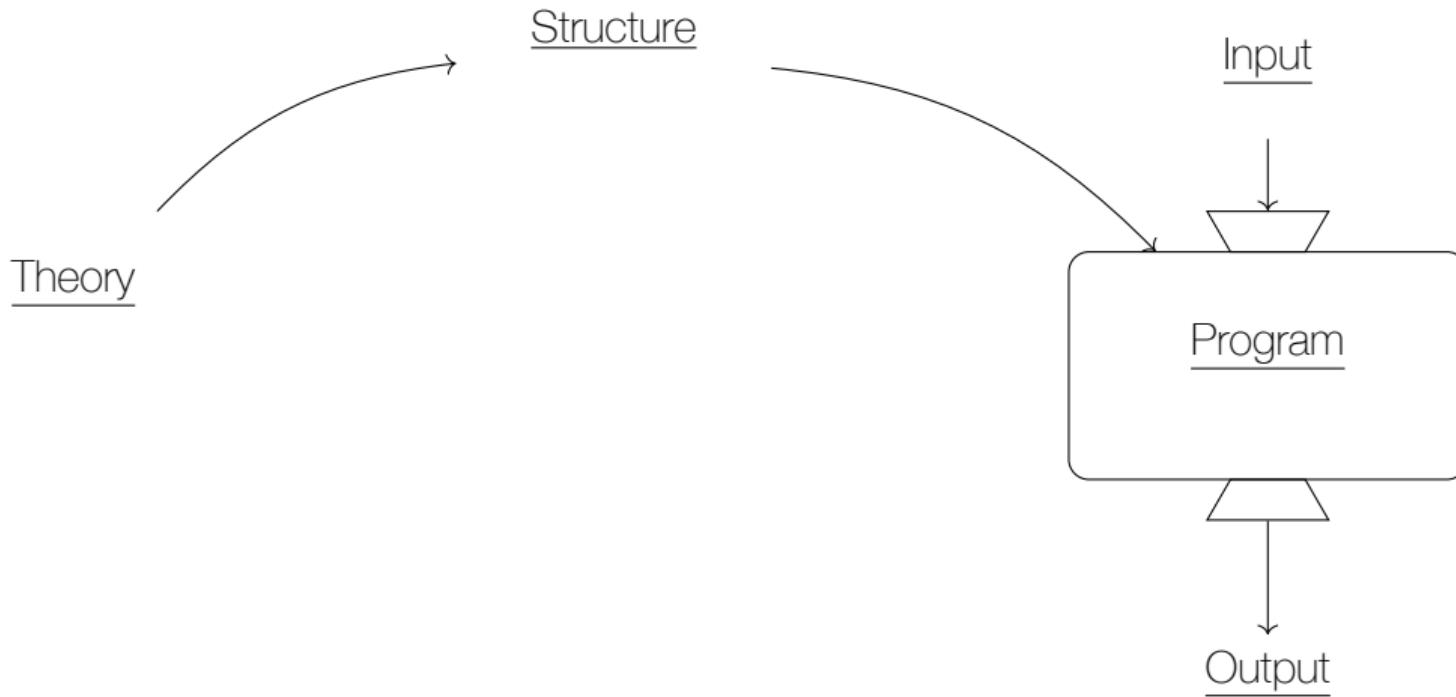
$$P := \lambda m. \lambda n. \lambda f. \lambda x. m f (n f x)$$

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

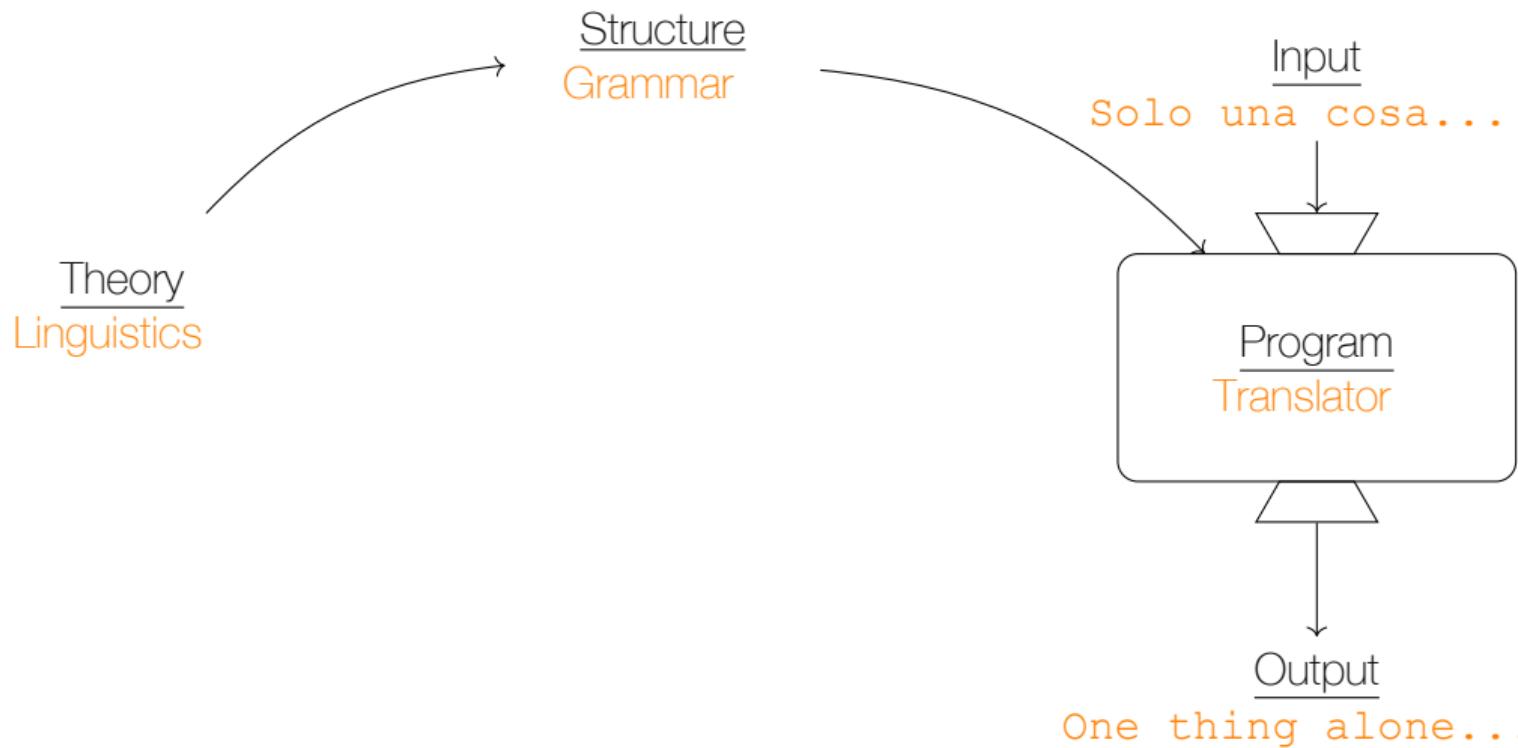
# Making It Explicit



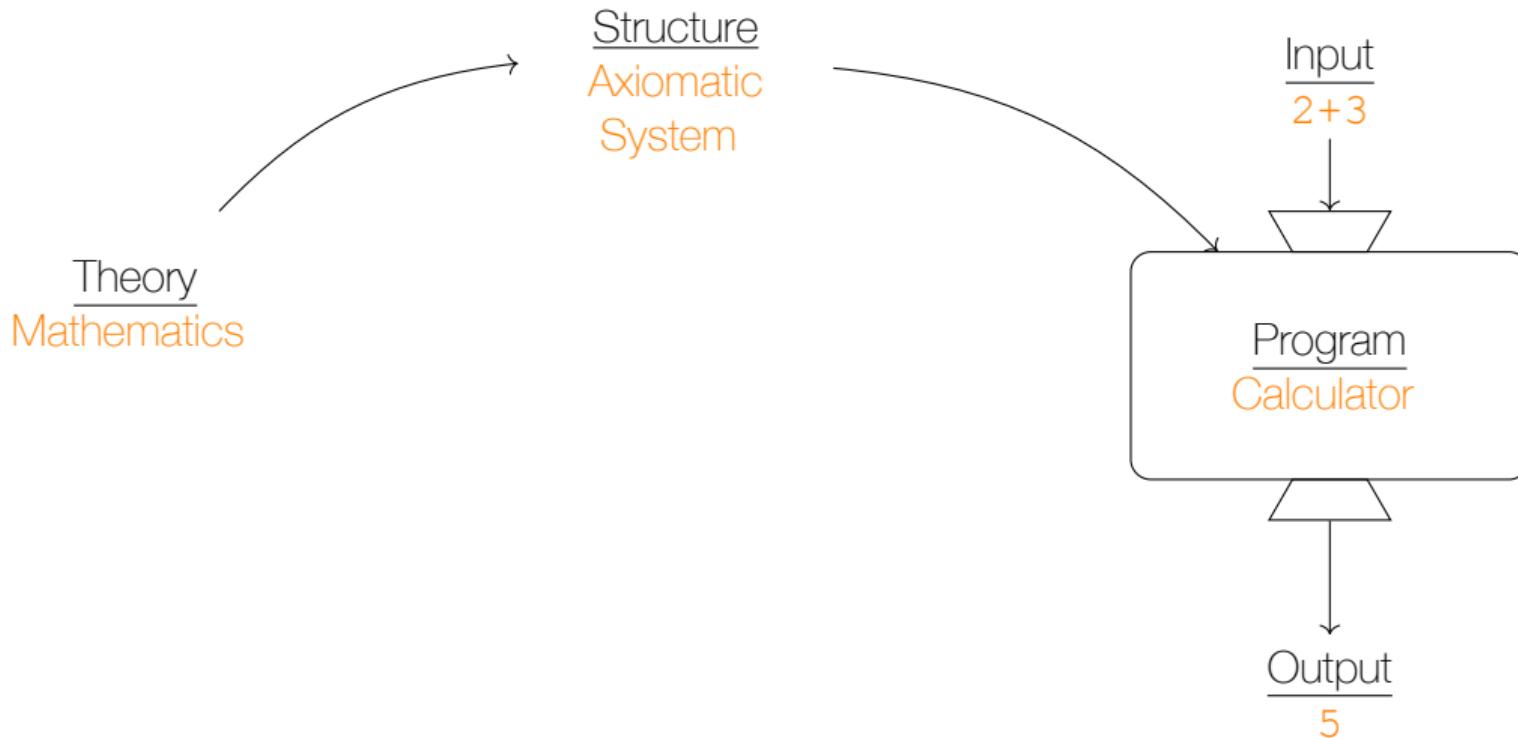
# Making It Explicit



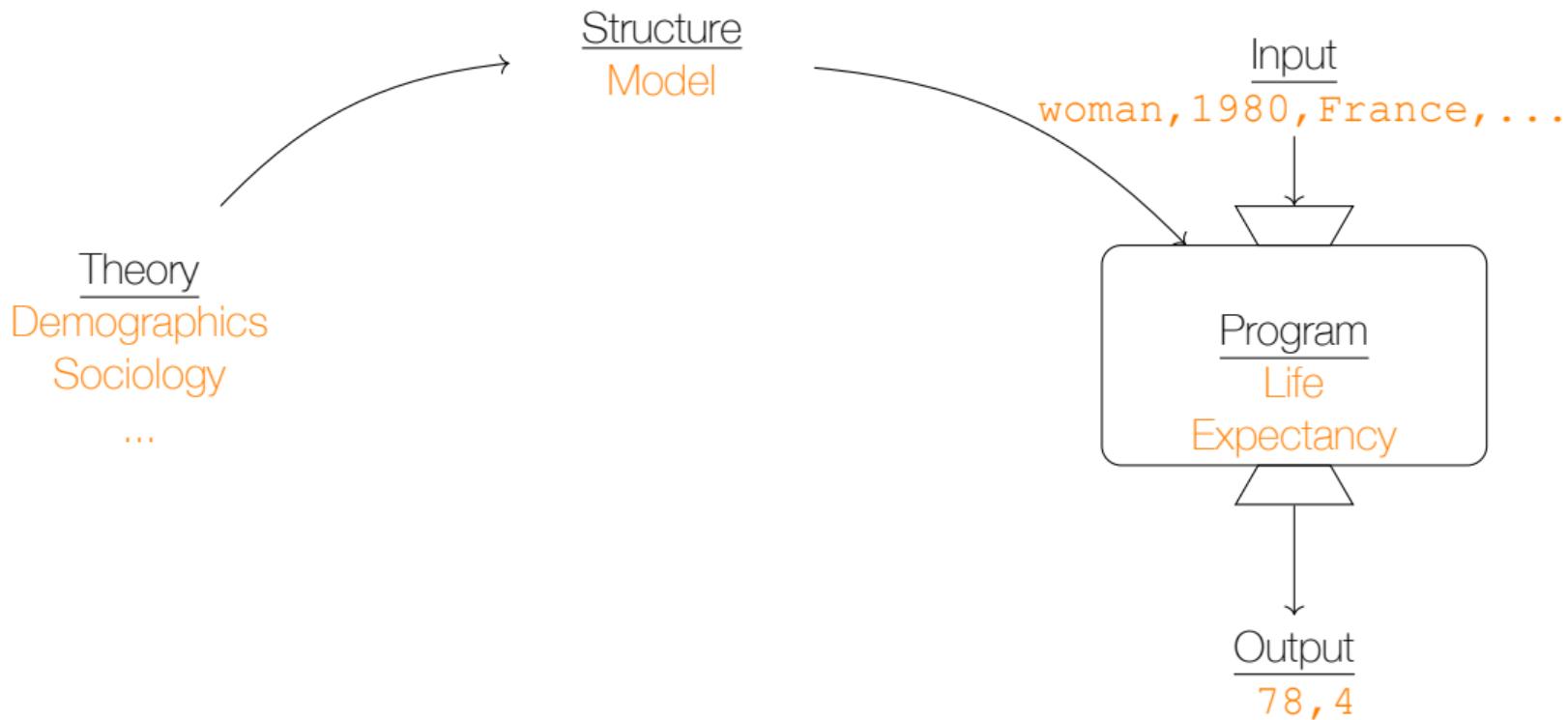
# Making It Explicit



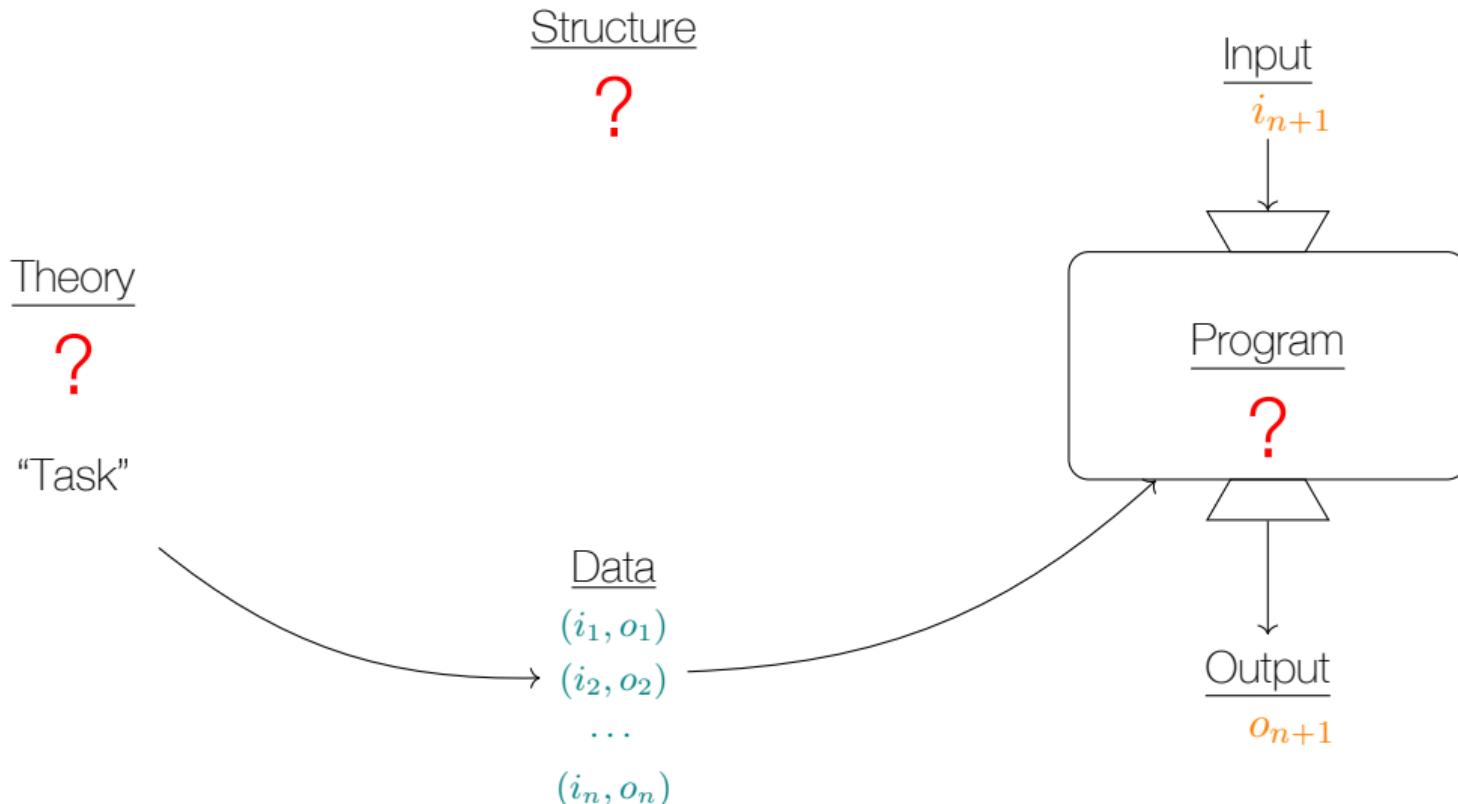
# Making It Explicit



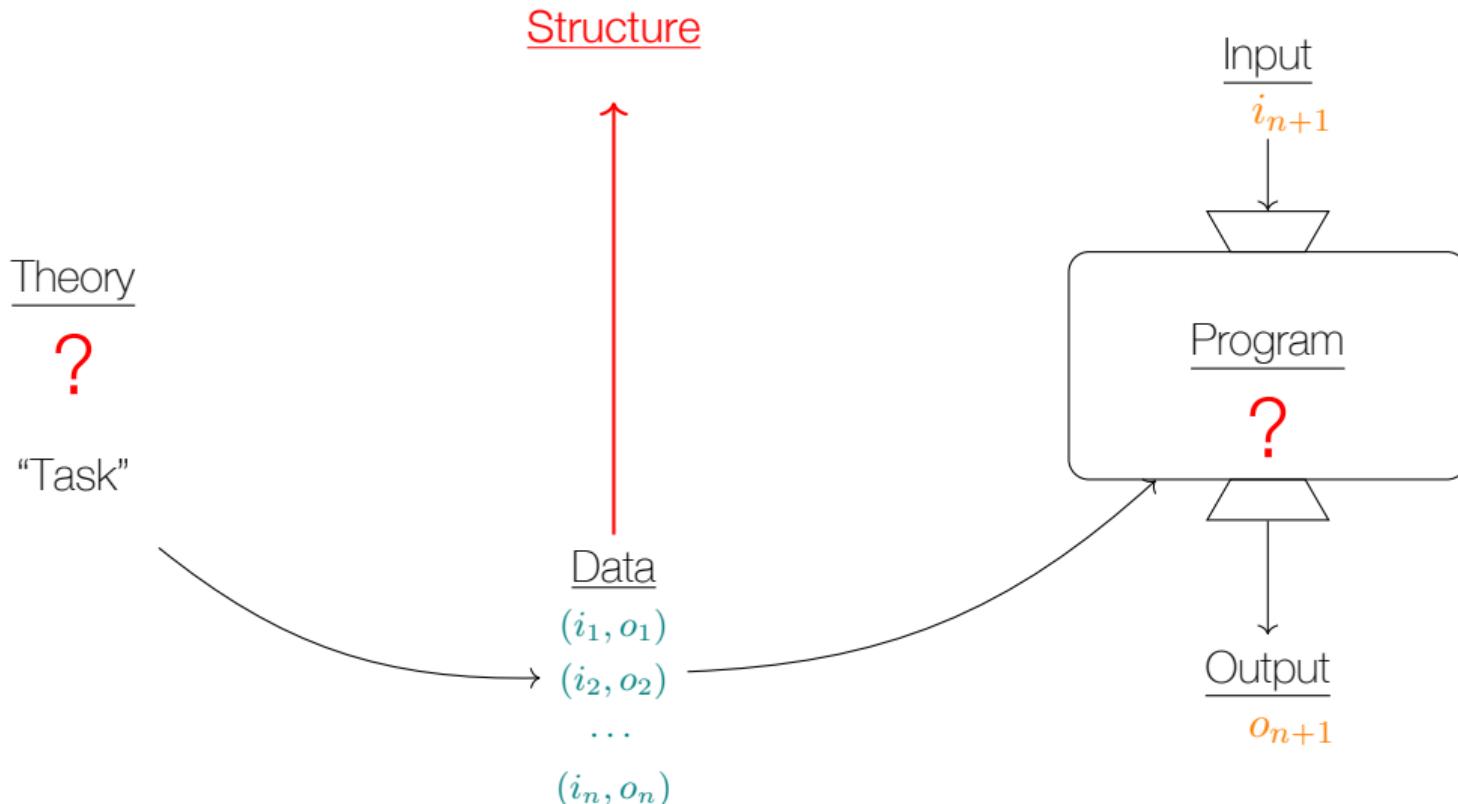
# Making It Explicit



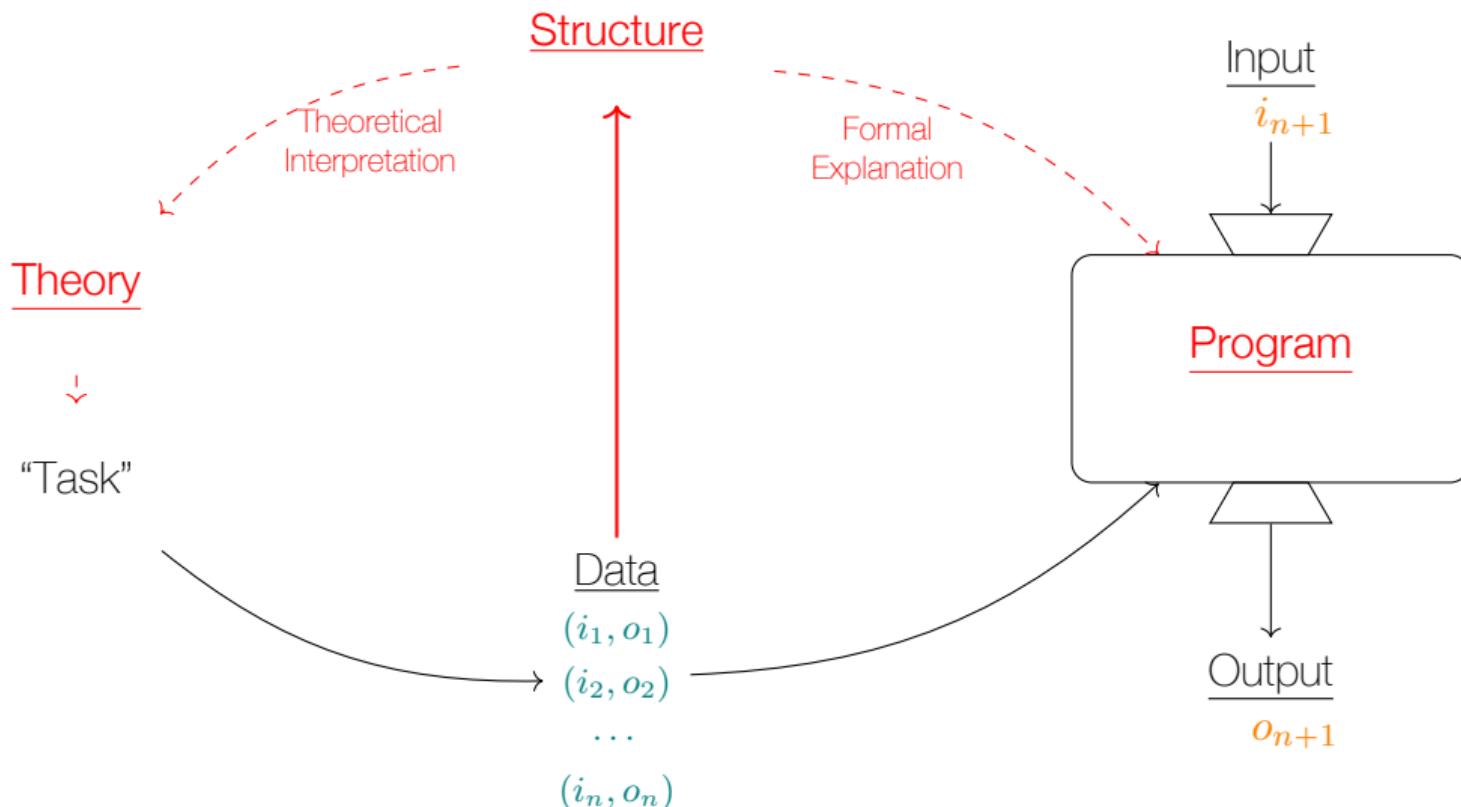
# Making It Explicit



# 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

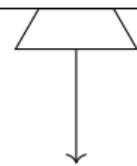
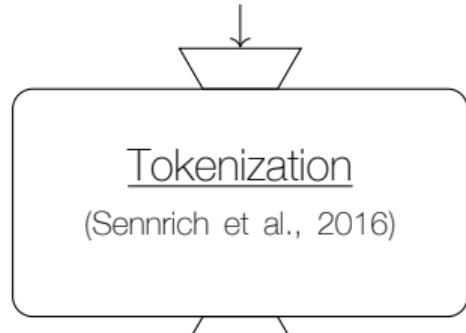
The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

# Formal Explainability

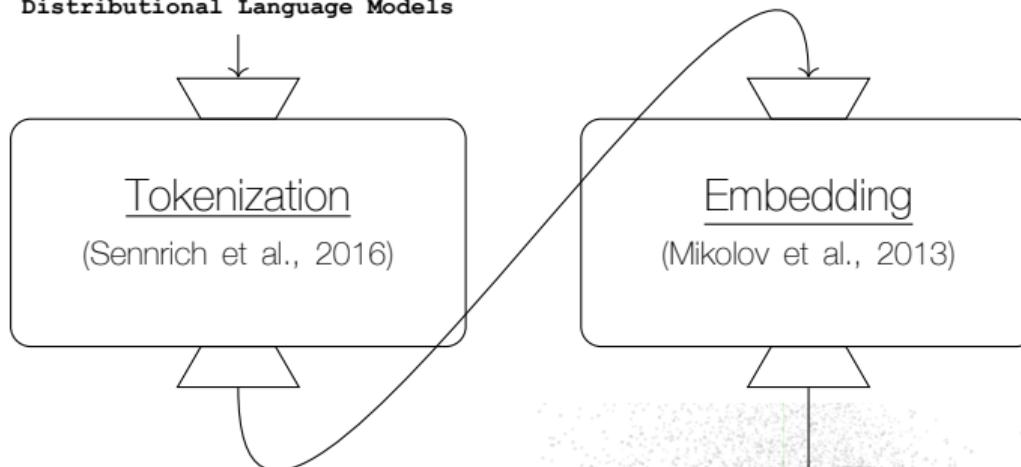
Epistemology of Machine Learning  
Distributional Language Models



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

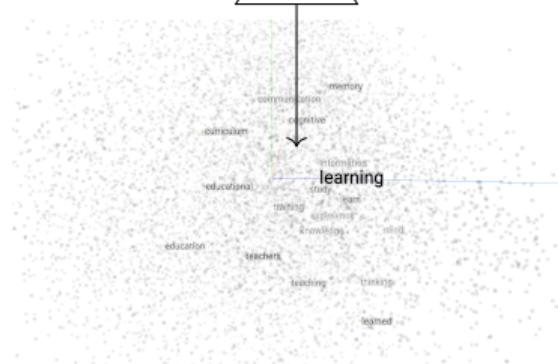
# Formal Explainability

Epistemology of Machine Learning  
Distributional Language Models



Epistemology of Machine Learning  
Distributional Language Models

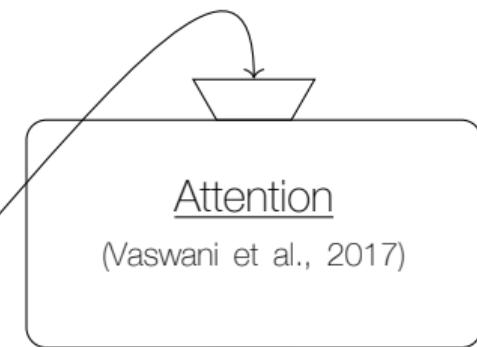
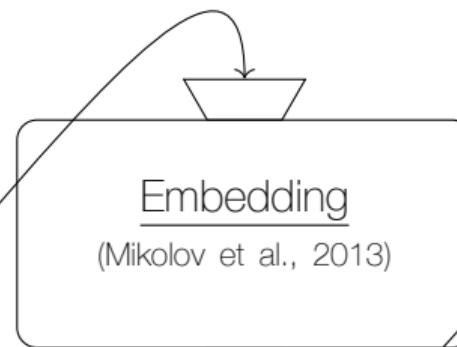
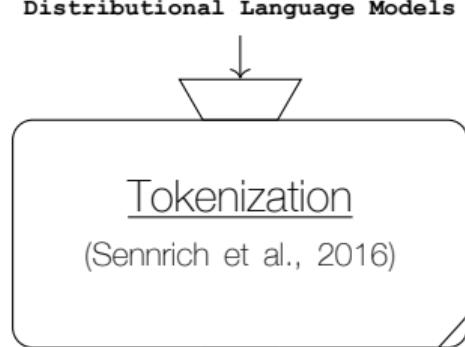
(<https://tiktoktokenizer.vercel.app>)



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

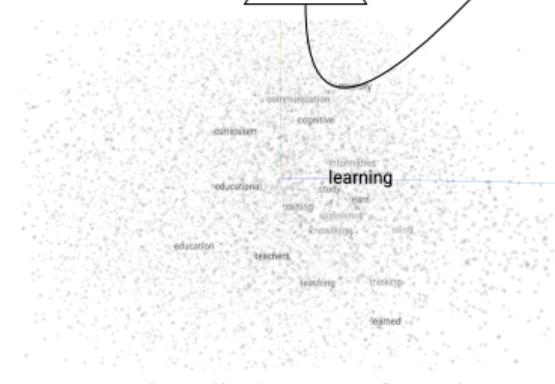
# Formal Explainability

Epistemology of Machine Learning  
Distributional Language Models



Epistemology of Machine Learning  
Distributional Language Models

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



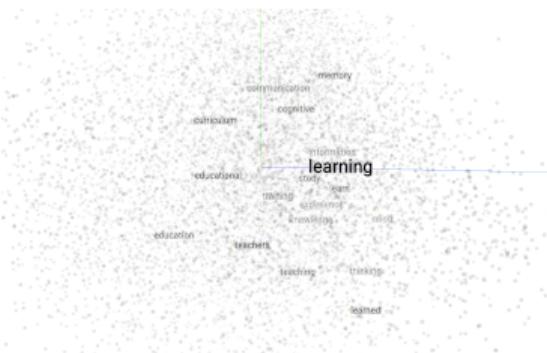
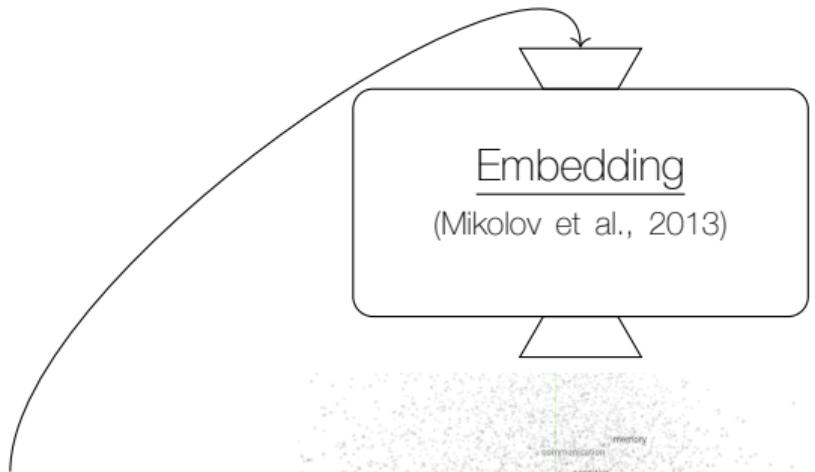
Ep  
ist  
em  
ology  
of  
Machine  
Learn  
ing  
Distribution  
al  
Language  
Models  
(<https://github.com/jessevieg/bertviz>)

A visualization showing the distribution of words in a sentence. The words "Epistemology of Machine Learning Distributional Language Models" are shown in a grid where each word has a unique color. The URL <https://github.com/jessevieg/bertviz> is displayed below the visualization.

# Formal Explainability

Epistemology of Machine Learning  
Distributional Language Models

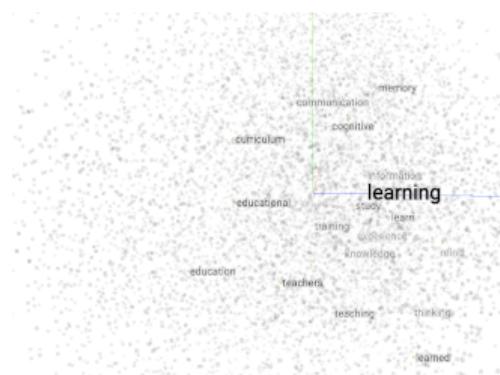
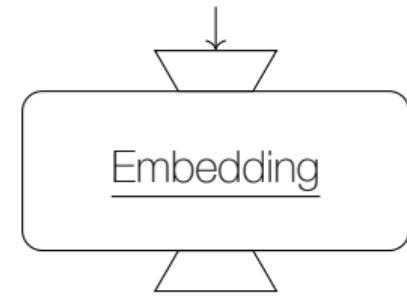
(<https://tiktoktokenizer.vercel.app>)



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

# The Structure of Embeddings

Epistemology of Machine Learning  
Distributional Language Models



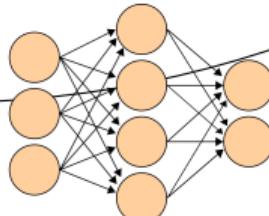
# The Structure of Embeddings

## Structure

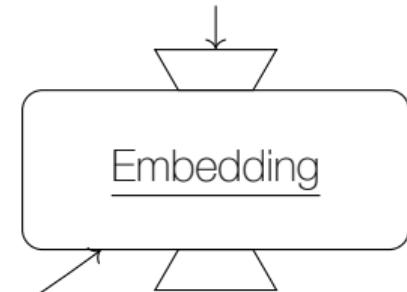
?

## Data

```
getElementsByClassName=function()
{
    if(t){(function(n,i,a){var g=n.getAttribute("id");
    (a,d,a)=c;});if(b.nodeType==9){return a;
    }on(a){return function(){var l=type(a);
    if(l=="text")return a.textContent;
    else return a.innerHTML;
    }}}
    return void 0==b||(ca.call(a,b)),type.function;
    if(e==b,apply(a[d],c),!1==breakValue(d));
    for(var c=a[b.length-1],d=a.length-1;b[c]==",",2),e=function(){return a.apply(b[0],b[1]);
    }((f=b.getElementsByTagName("div").nodeType||1)[0].parentNode||0);
    ["for","(h1.length-1);h1.length-1,h1.length-1];
    "for(c=a.split(";"),c.length-1,c.length-1);
    "for(var d,f=a[1],c.length-1,d.length-1,f.length-1;
    on(c,s){for(var d,f=a[1],c.length-1,d.length-1;
    nodeType|0)(v=b[0][i]);if(v.nodeType==1){for(var s=j.node_type||1;j<1;s+=1);
    or(var h=f[1];h.length-1,h.length-1,h.length-1);
    "for(var v=a[n],da,d,v.length-1,v.length-1,v.length-1;
    v.length-1);return va(n),da,d,v.length-1,v.length-1;
    v.length-1);}else{var v=a[n],da,d,v.length-1,v.length-1;
    v.length-1);}});}}
```



Epistemology of Machine Learning  
Distributional Language Models



Embedding

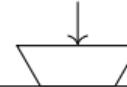


# The Structure of Embeddings

Structure

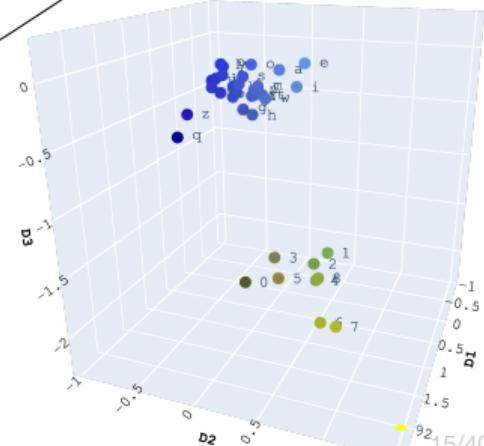
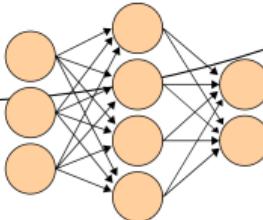
?

{-, /, 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)

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

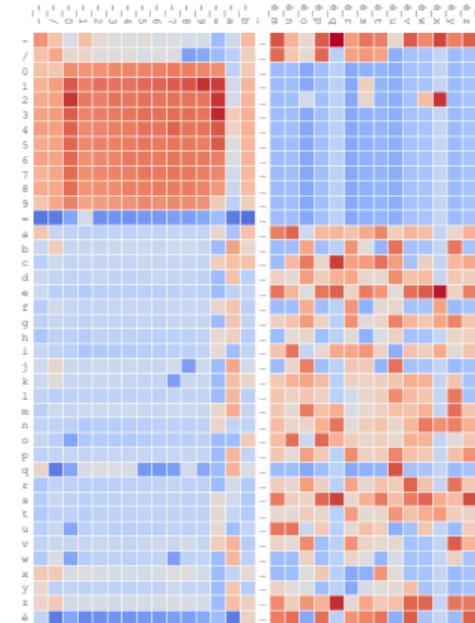
$$\frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{when} \quad \vec{w} \cdot \vec{c} = \log \left( \frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

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

# Example: Characters in Wikipedia

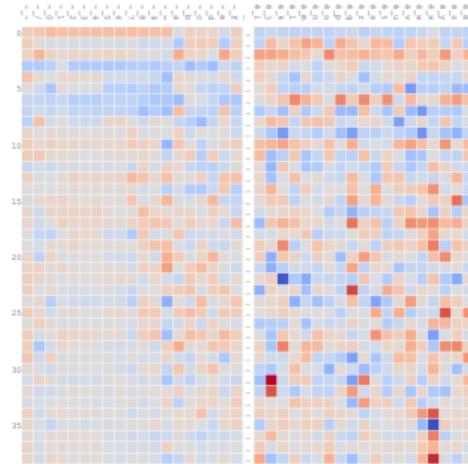
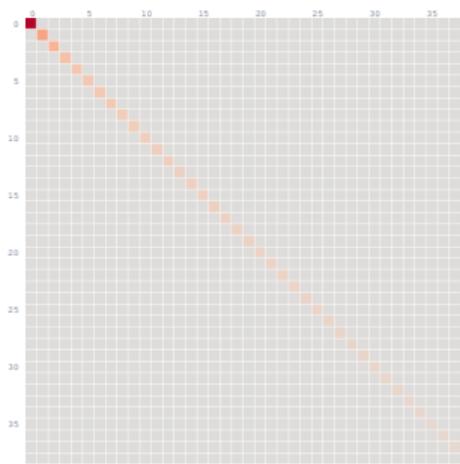
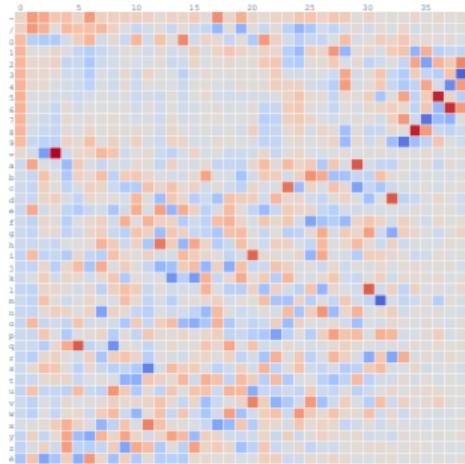
$$W = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, é\}$$

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



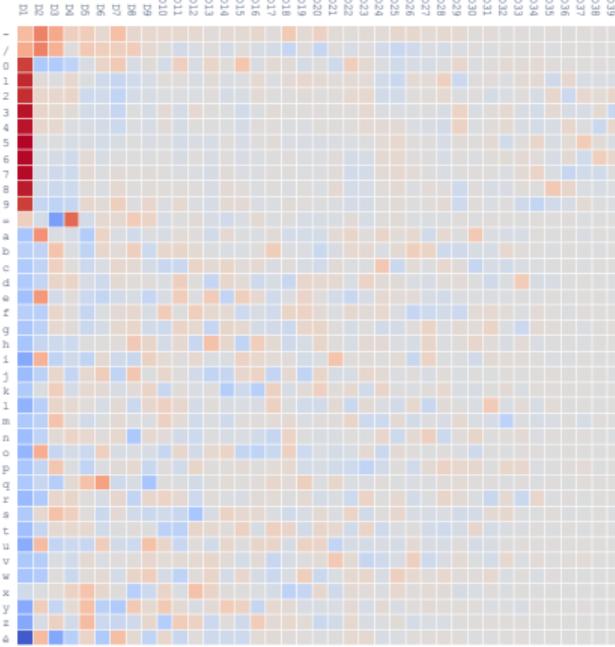
$$\begin{aligned} M_{wc} &= \text{pmi}(w, c) \\ &= \log \frac{p(w, c)}{p(w)p(c)} \end{aligned}$$

# SVD of Wikipedia Character PMI Matrix

 $U$  $\Sigma$  $V^T$ 

Truncate

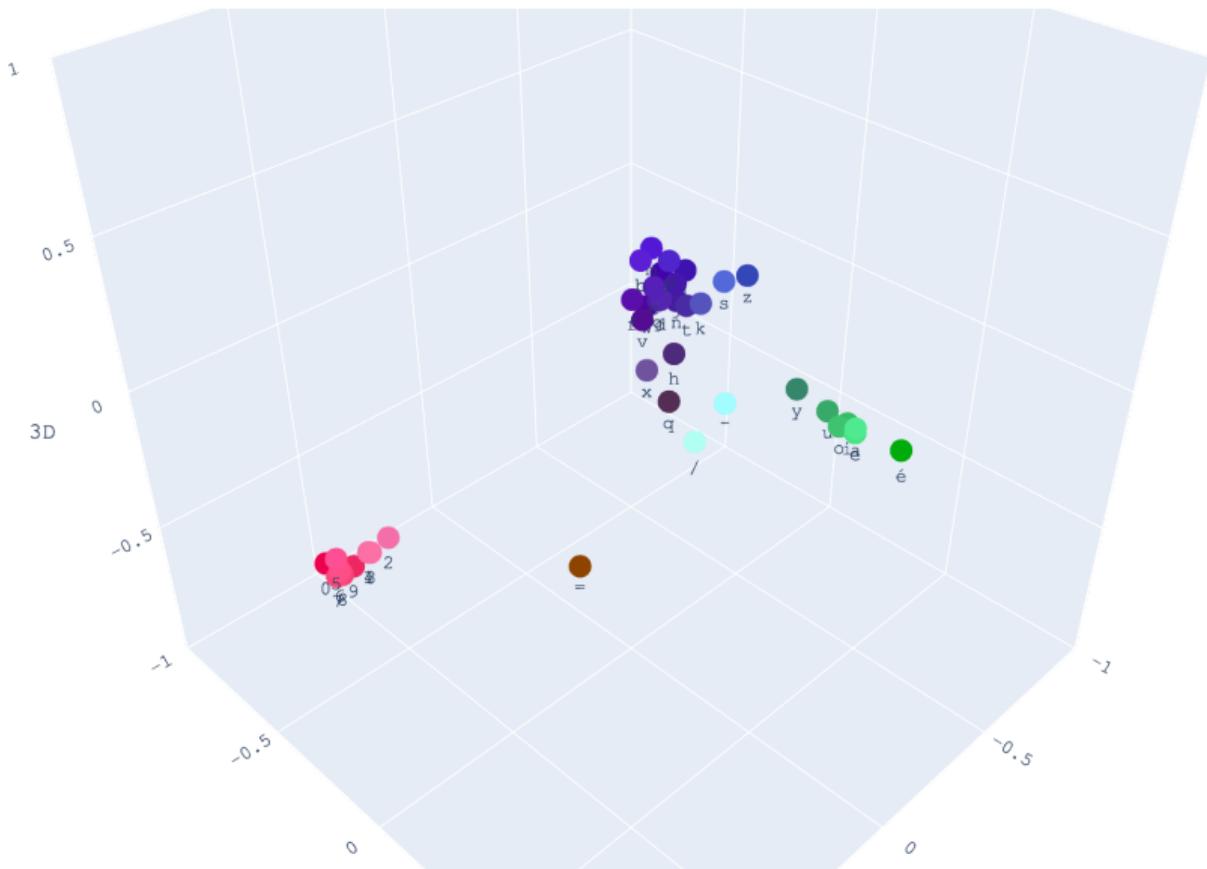
$U \times \Sigma$



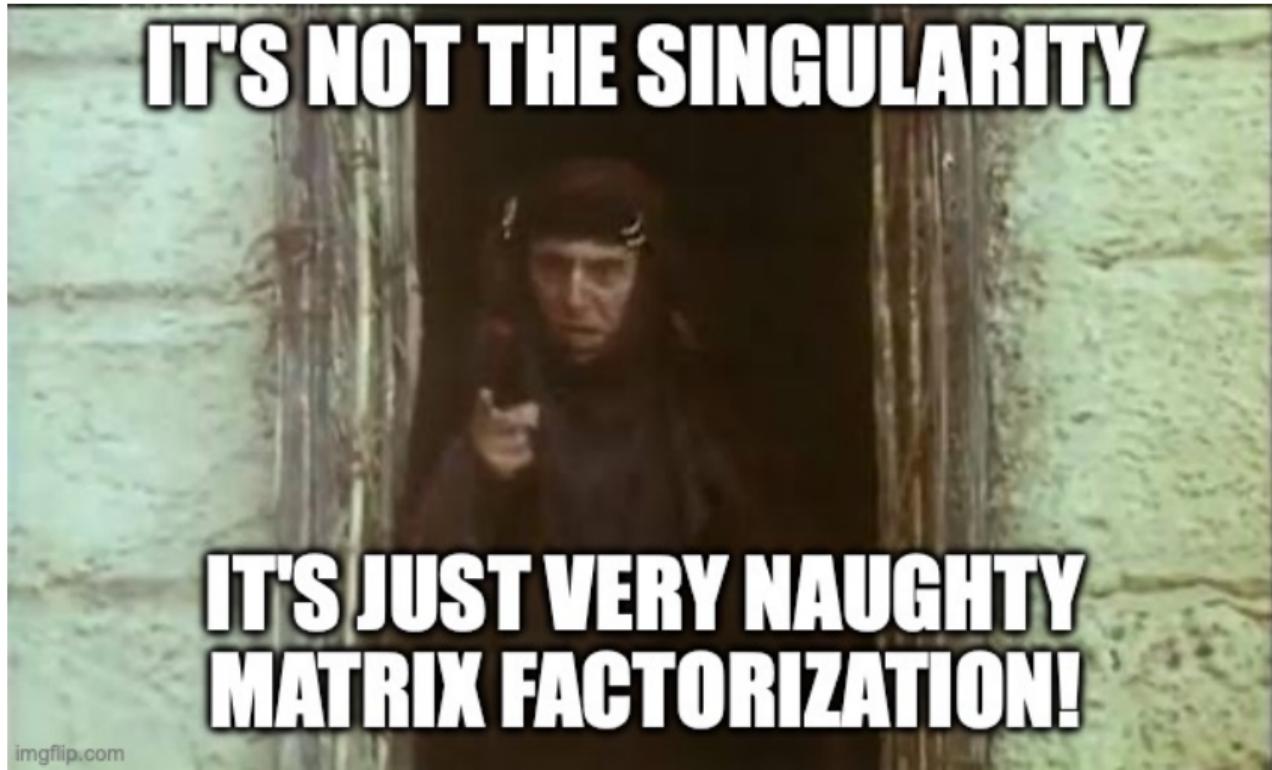
Truncate

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



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

What to conclude?



## 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.

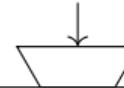
(Goldberg and Levy, 2014)

# The Structure of Embeddings

Structure

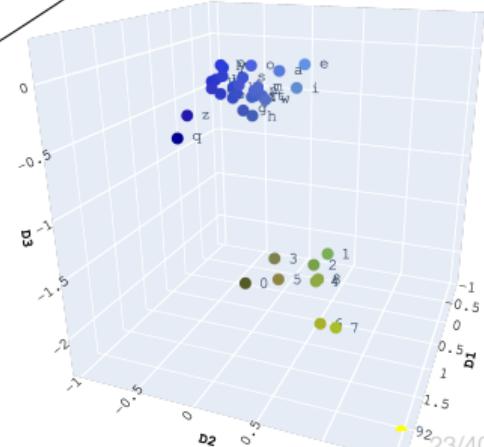
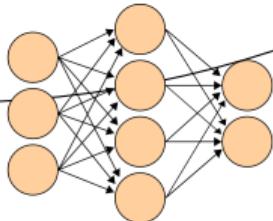
?

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



Embedding

Data



# The Structure of Embeddings

## Structure

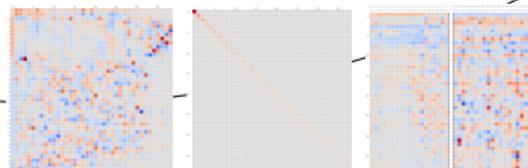


$\{-, /, 0, 1, 2, \dots, 8, 9, =,$   
 $a, b, c, \dots, w, x, y, z, \acute{e}\}$

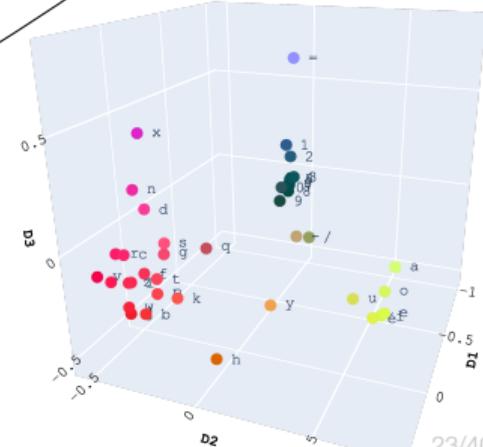


Embedding

## Data



SVD



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

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

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

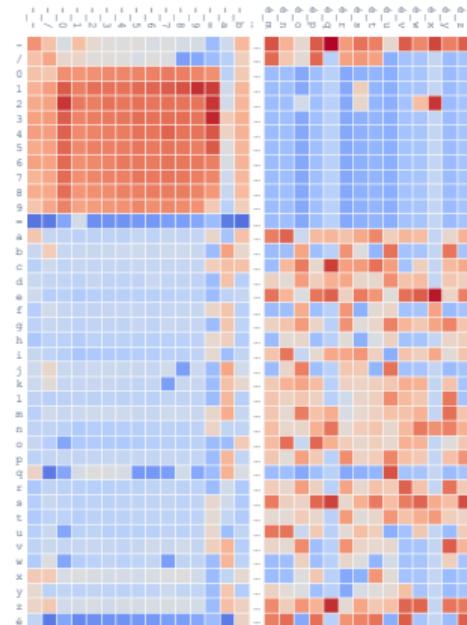
# Embeddings as Functions Over Sets

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

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

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

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



# Embeddings as Functions Over Sets

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

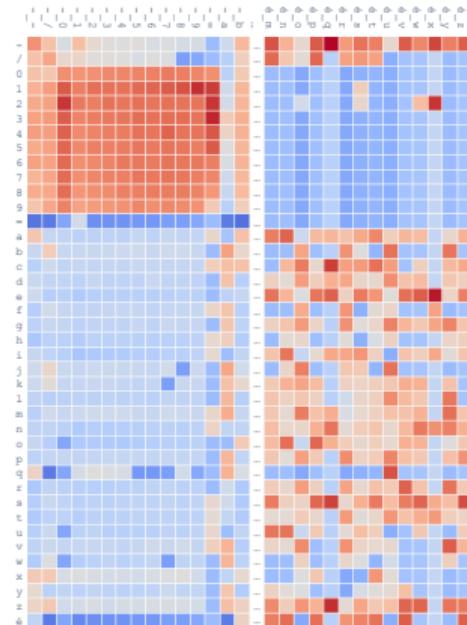
$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, z), (\text{é}, \text{é})\}$$

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

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

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

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



# Embeddings as Functions Over Sets

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

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

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

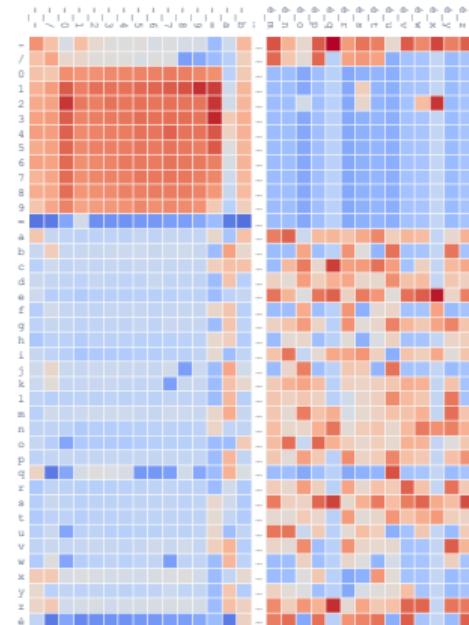
$$(\textcolor{red}{x}, \textcolor{blue}{y}) \mapsto \text{pmi}(\textcolor{red}{x}, \textcolor{blue}{y})$$

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

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

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

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



# 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

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

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

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

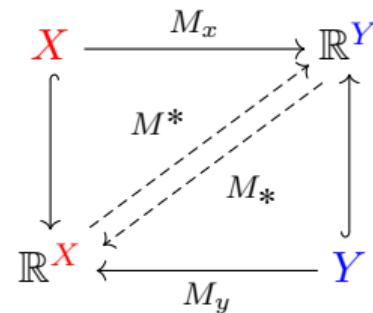
$$(\textcolor{red}{x}, \textcolor{blue}{y}) \mapsto \text{pmi}(\textcolor{red}{x}, \textcolor{blue}{y})$$

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

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

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

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

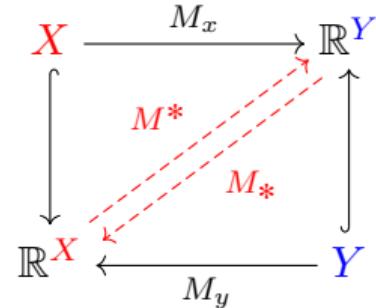


$$M^*: \mathbb{R}^{\textcolor{red}{X}} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

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

# Embeddings as Functions Over Sets

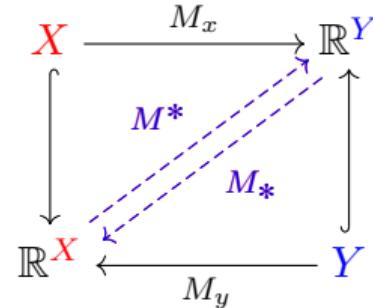
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$



# 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

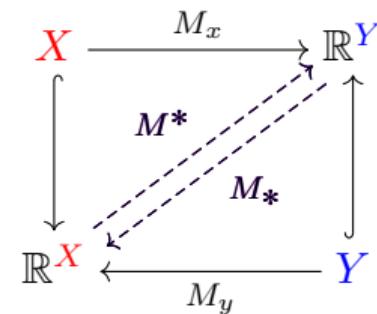
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

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

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

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

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



# Embeddings as Functions Over Sets

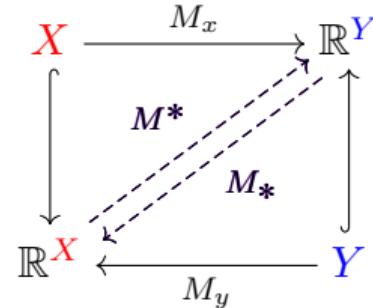
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

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

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

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

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



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

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

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

# Embeddings as Functions Over Sets

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

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

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

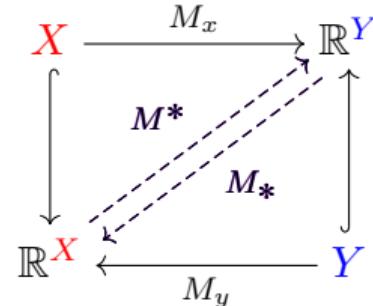
$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

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

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

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

The  $u_i$  and  $v_i$  are (linear)  
fixed points!



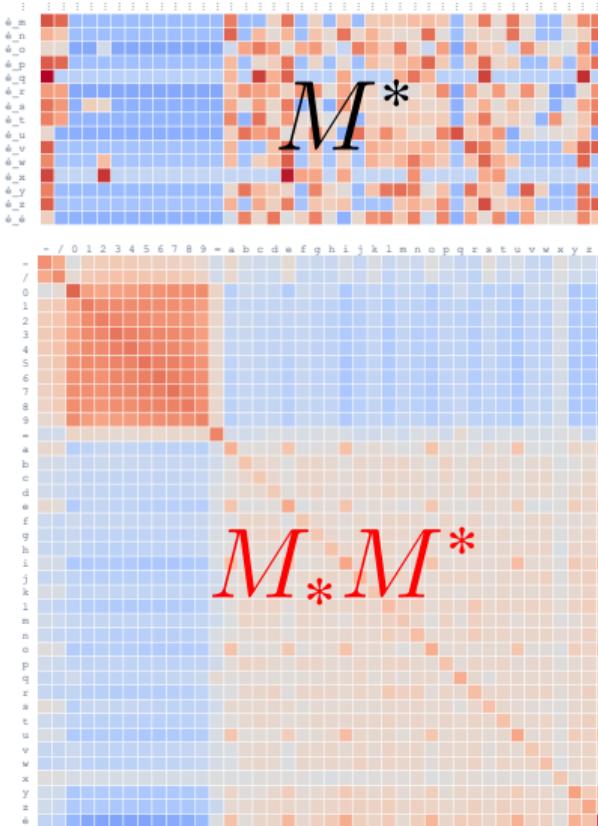
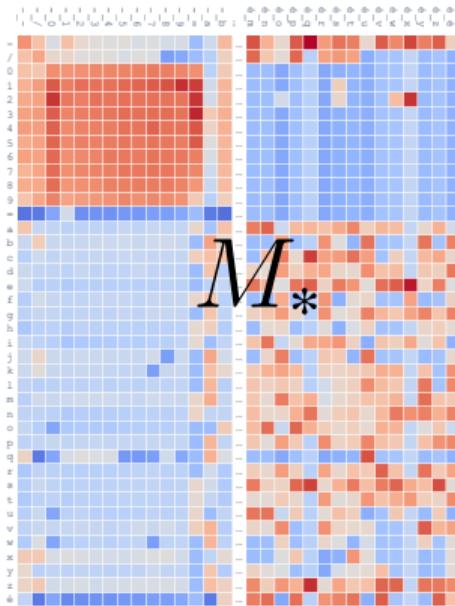
$$U := [u_1, \dots, u_m]$$

$$V := [v_1, \dots, v_n]$$

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

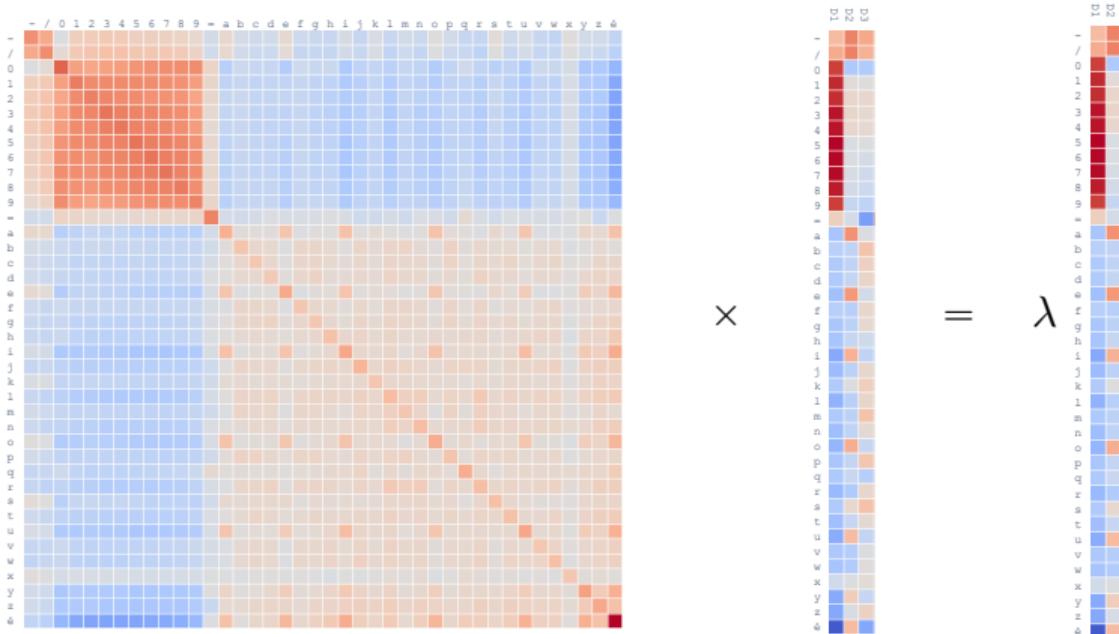
$$M = U \Sigma V^T$$

# $M_* M^*$ as a Covariance Matrix

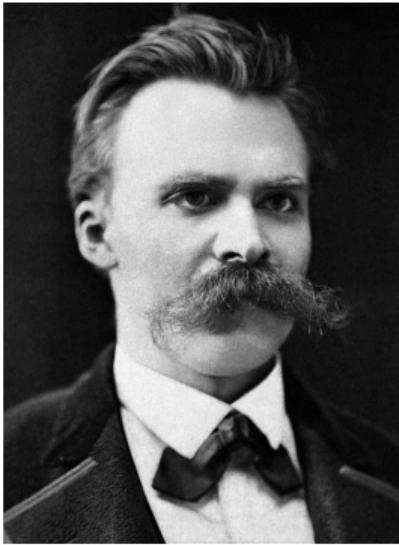


# Eigenvectors as Fixed Points

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



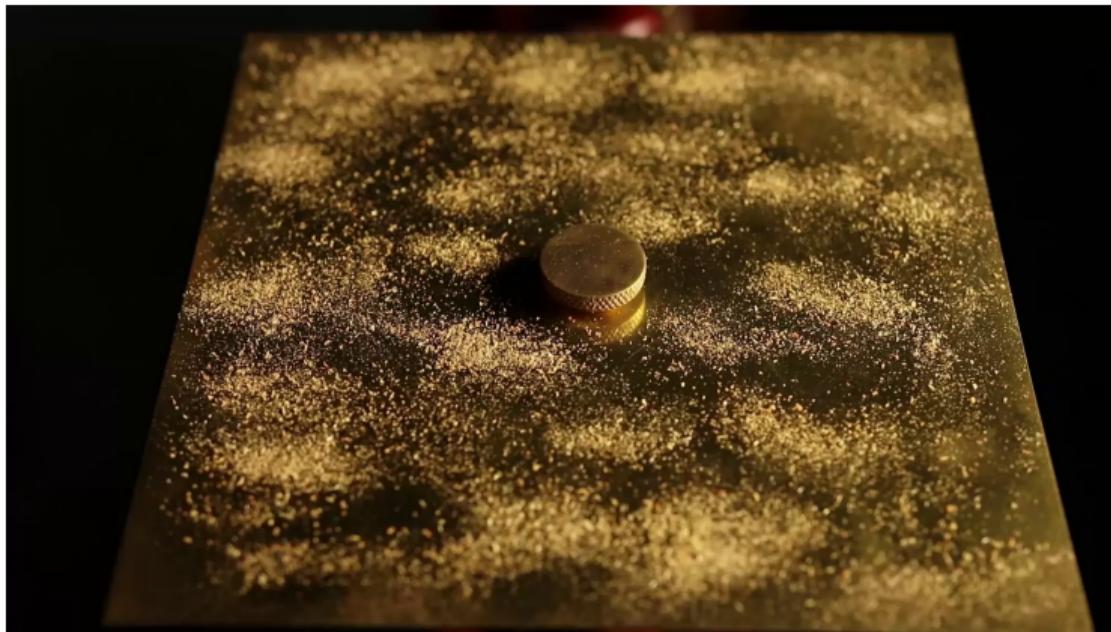
# Chladni Figures



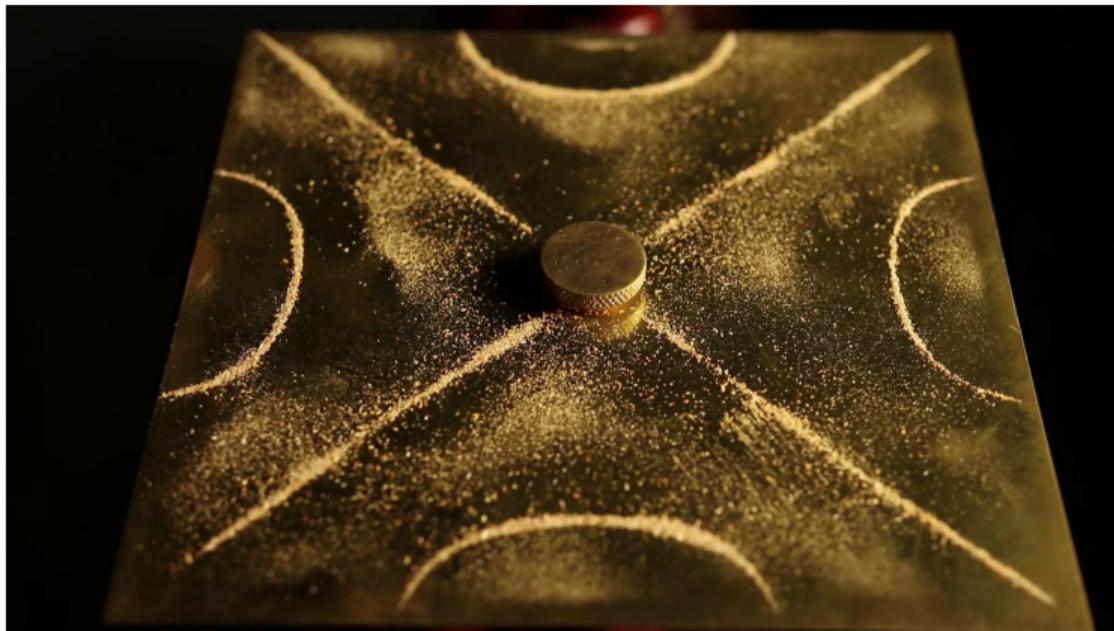
“One can conceive of a profoundly **deaf** human being who has never experienced sound or music; just as such a person will gaze in astonishment at the **Chladnian sound-figures in sand**, find their cause in the vibration of a string, and swear that **he must now know what men call sound** — this is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

(Nietzsche, 1873)

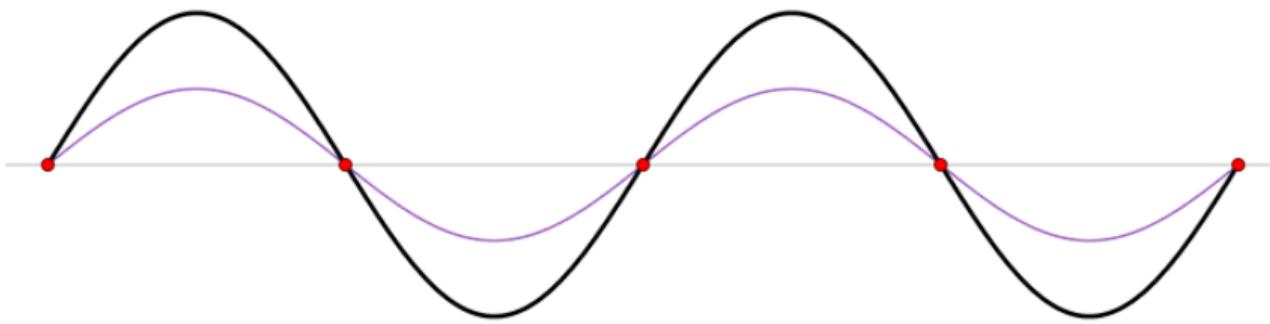
# Chladni Figures



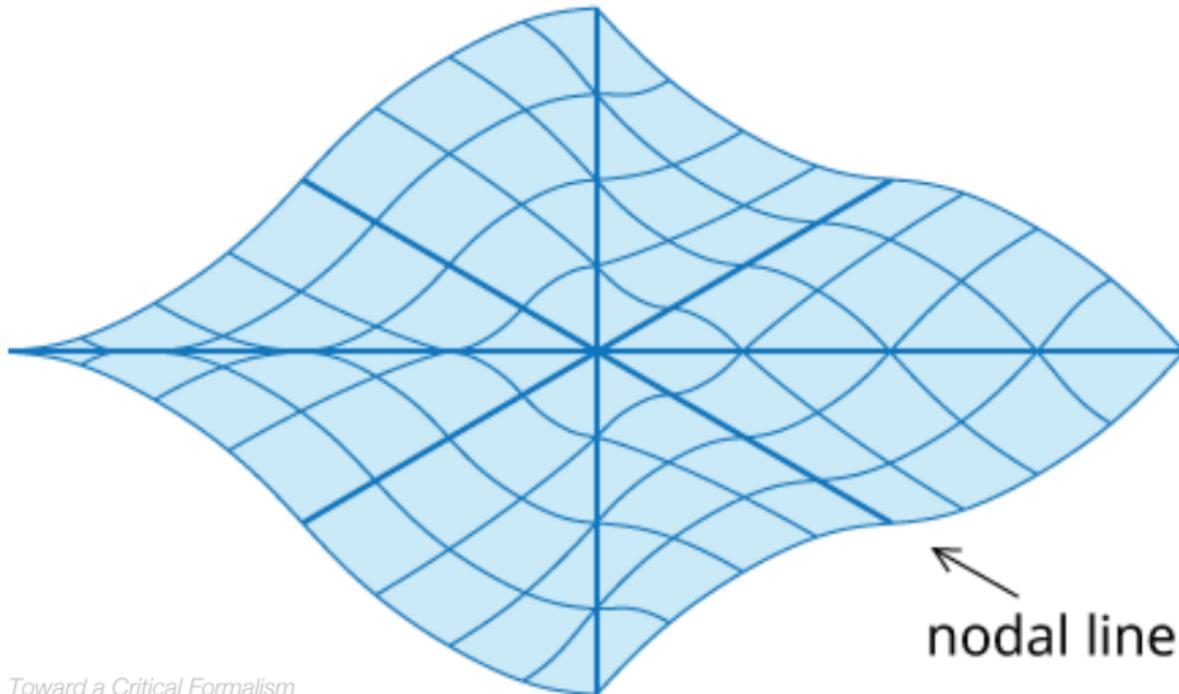
# Chladni Figures



# Chladni Figures



## Chladni Plate

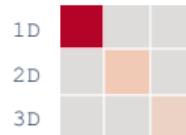


# Structural Features

Eigenvectors of  $M_* M^*$ :



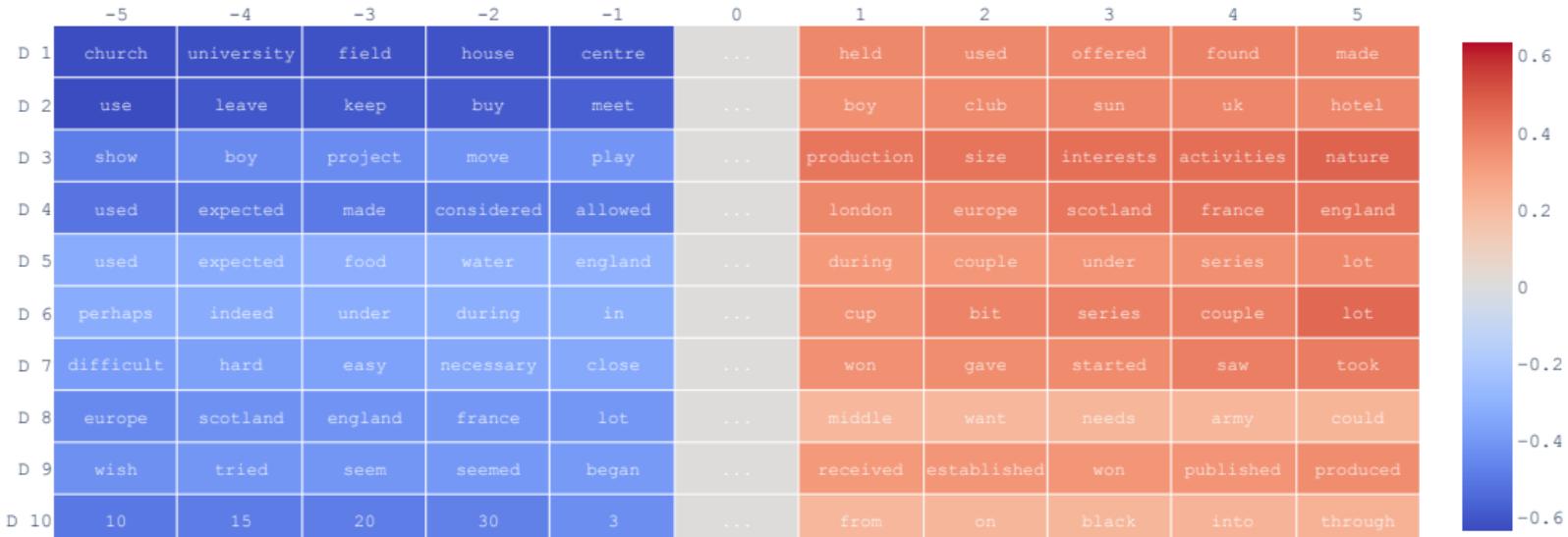
Eigenvalues of  $M_* M^*$  and  $M^* M_*$ :



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

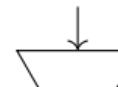
Conclusion

# The Structure of Embeddings

## Structure

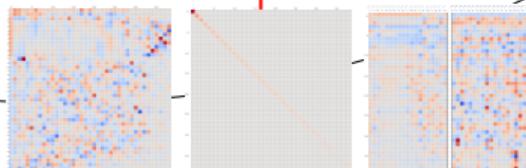


$\{-, /, 0, 1, 2, \dots, 8, 9, =,$   
 $a, b, c, \dots, w, x, y, z, \acute{e}\}$

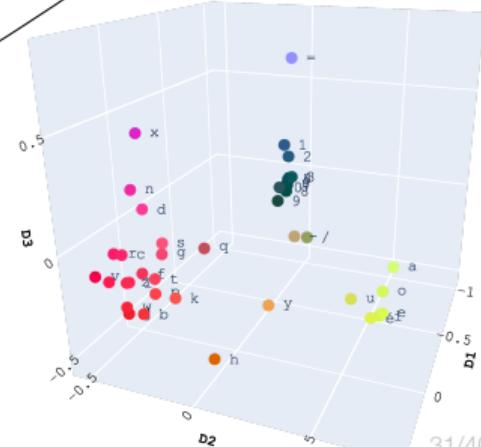


Embedding

## Data



SVD

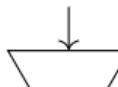


# The Structure of Embeddings

## Structure

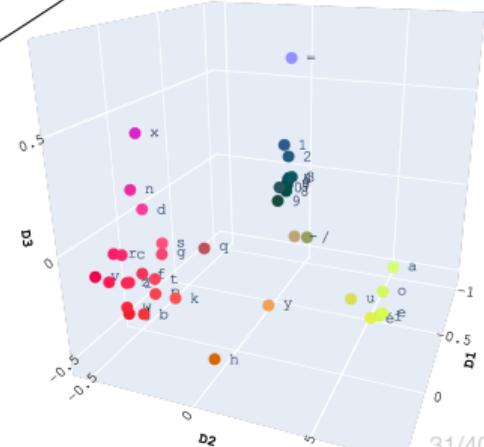
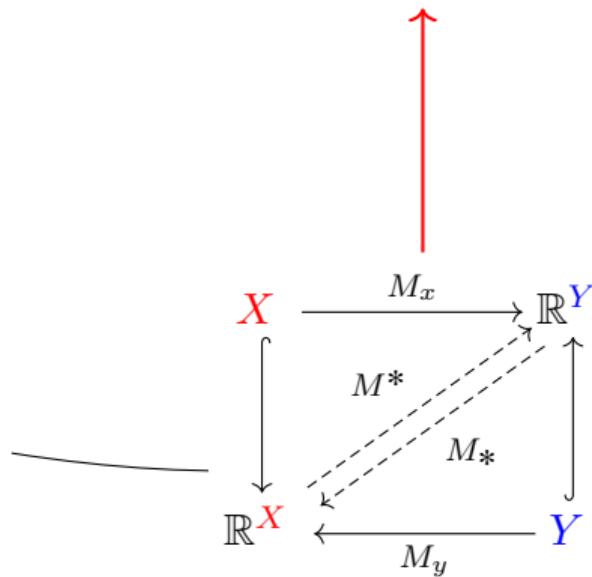


$\{-, /, 0, 1, 2, \dots, 8, 9, =,$   
 $a, b, c, \dots, w, x, y, z, \acute{e}\}$



Embedding

## Data



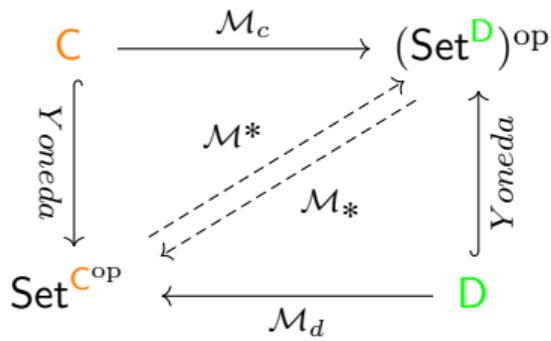
# The Structure of Embeddings

Structure

?

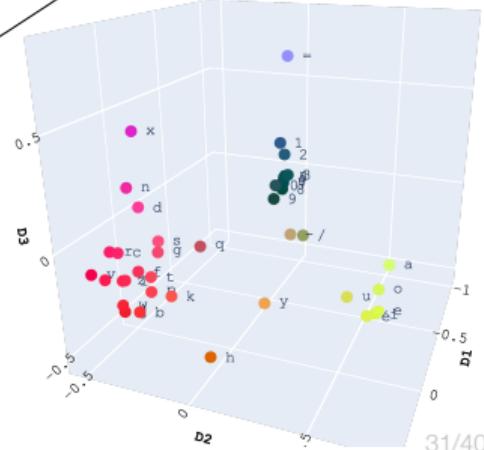
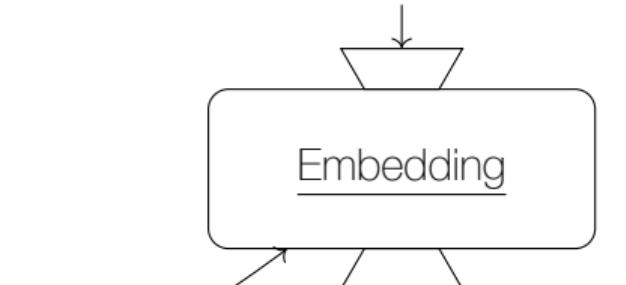


Data



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

Embedding

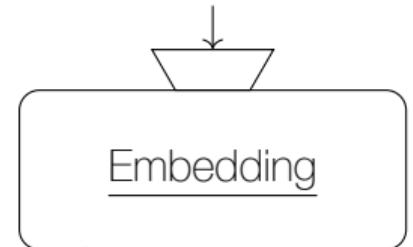


# The Structure of Embeddings

Structure

?

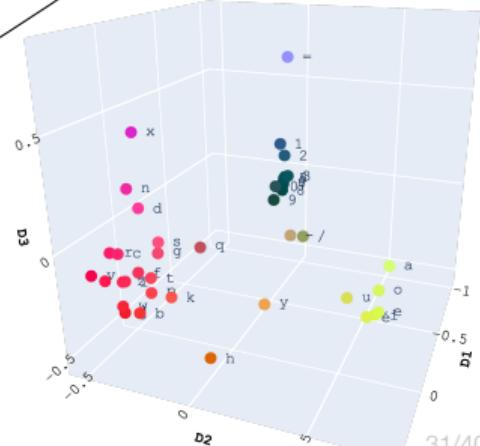
{-, /, 0, 1, 2, ..., 8, 9, =,  
a, b, c, ..., w, x, y, z, é}



Data



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



Structure

?

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

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

?

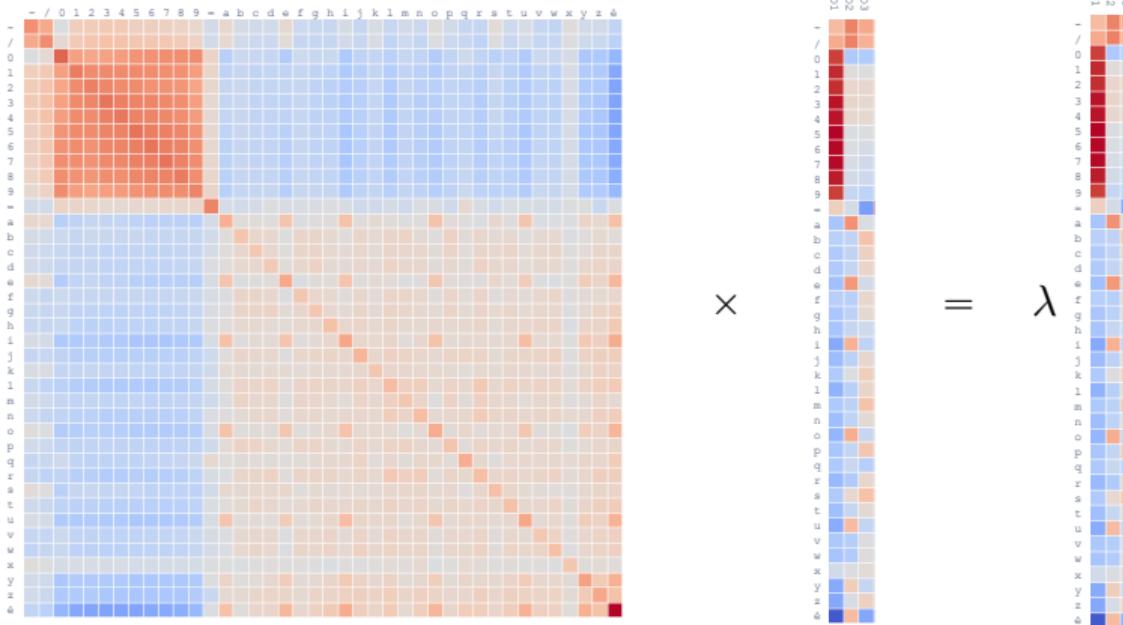
$$\begin{array}{ccc} \textcolor{teal}{term}_i & \textcolor{teal}{context}_i & \text{measure} \\ \downarrow & \downarrow & \swarrow \\ \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \textcolor{red}{2} \end{array}$$

Structure

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

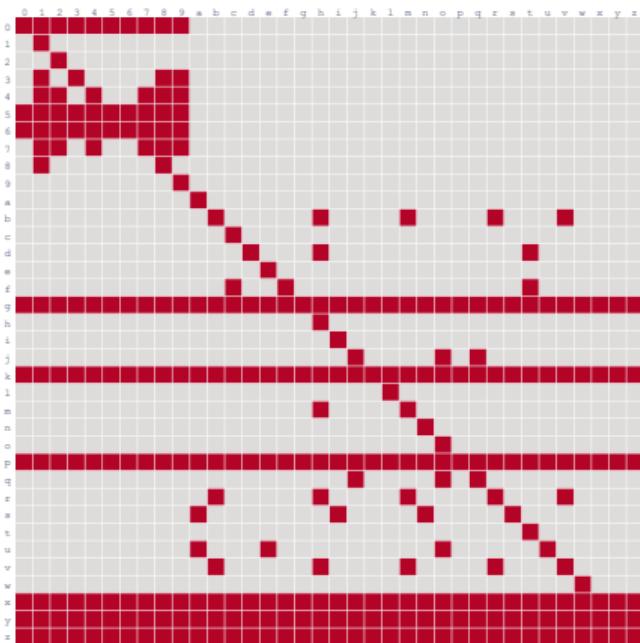
# Binary Fixed Points

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



# Binary Fixed Points

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



★

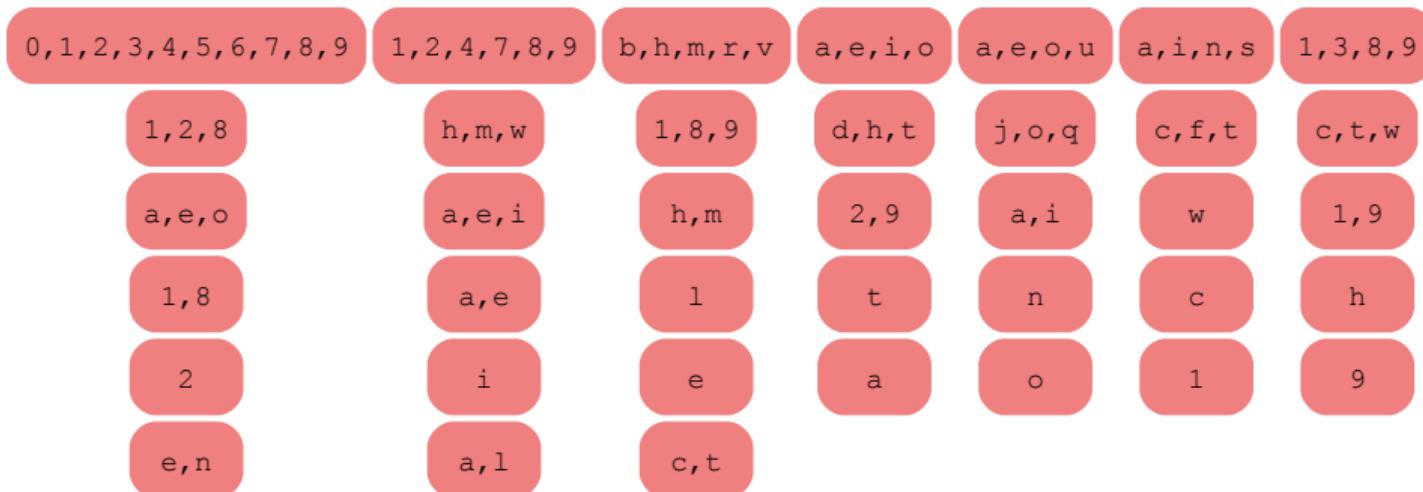


?

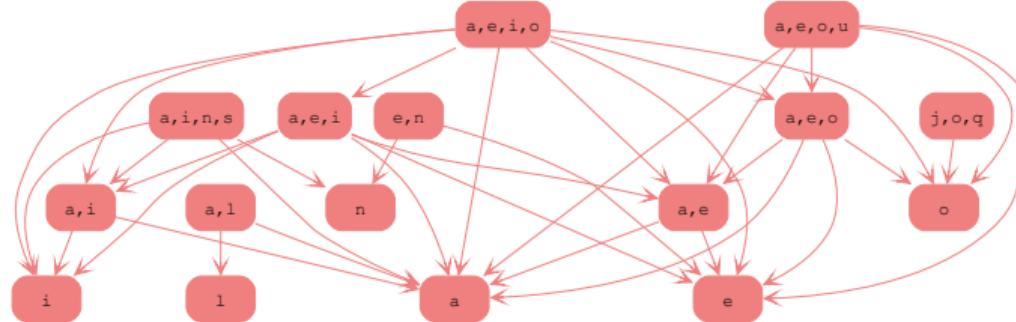
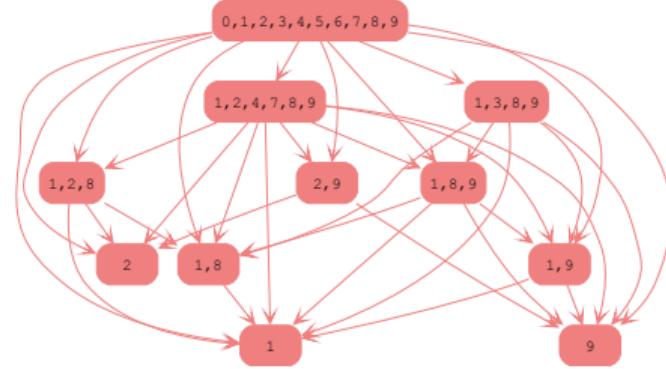
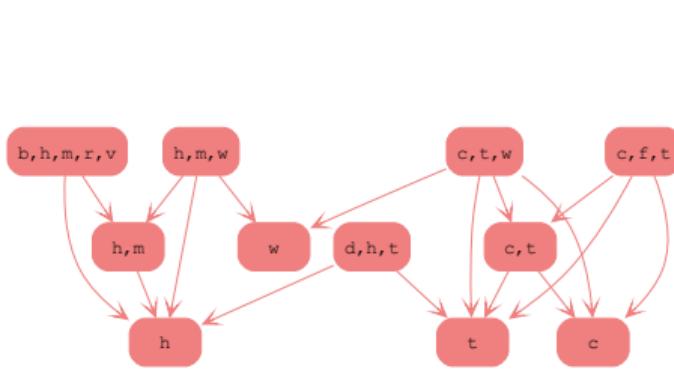


# “Eigensests”

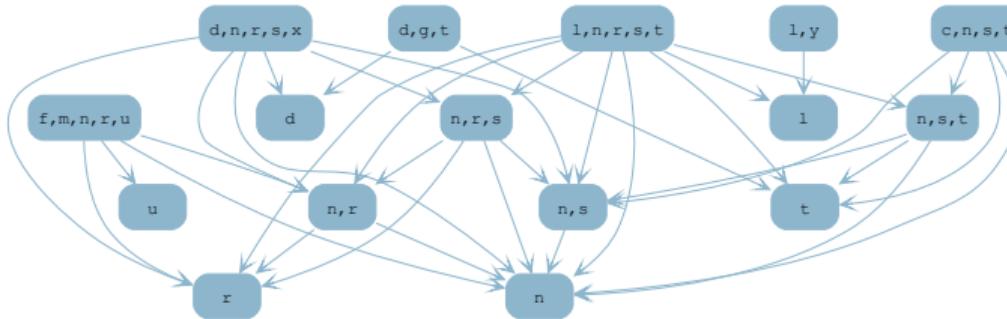
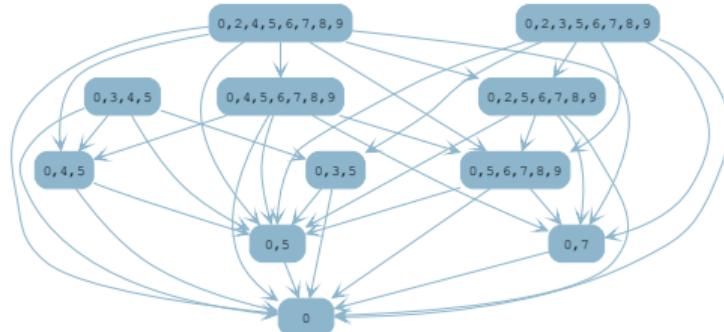
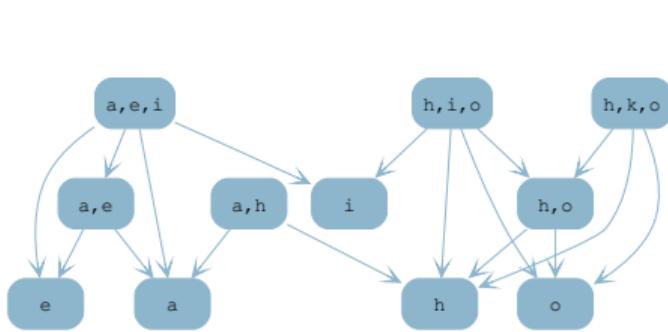
$$\mathcal{M}_*\mathcal{M}^*f = f$$



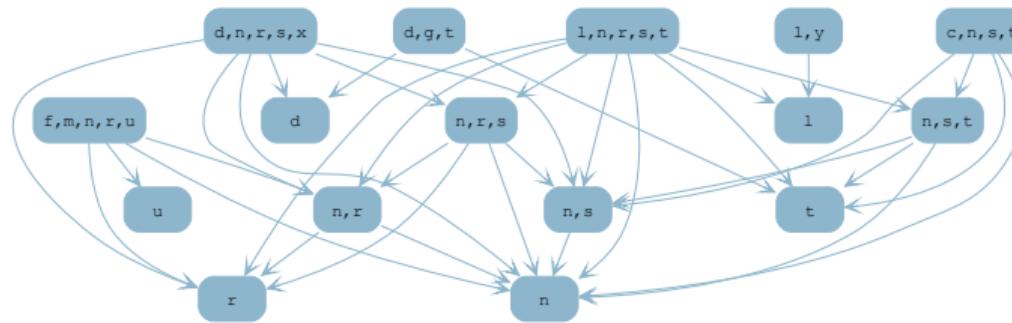
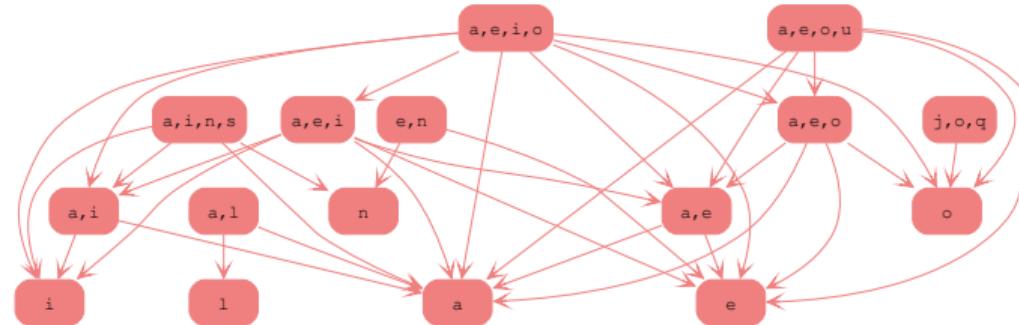
# Partial Order Structure

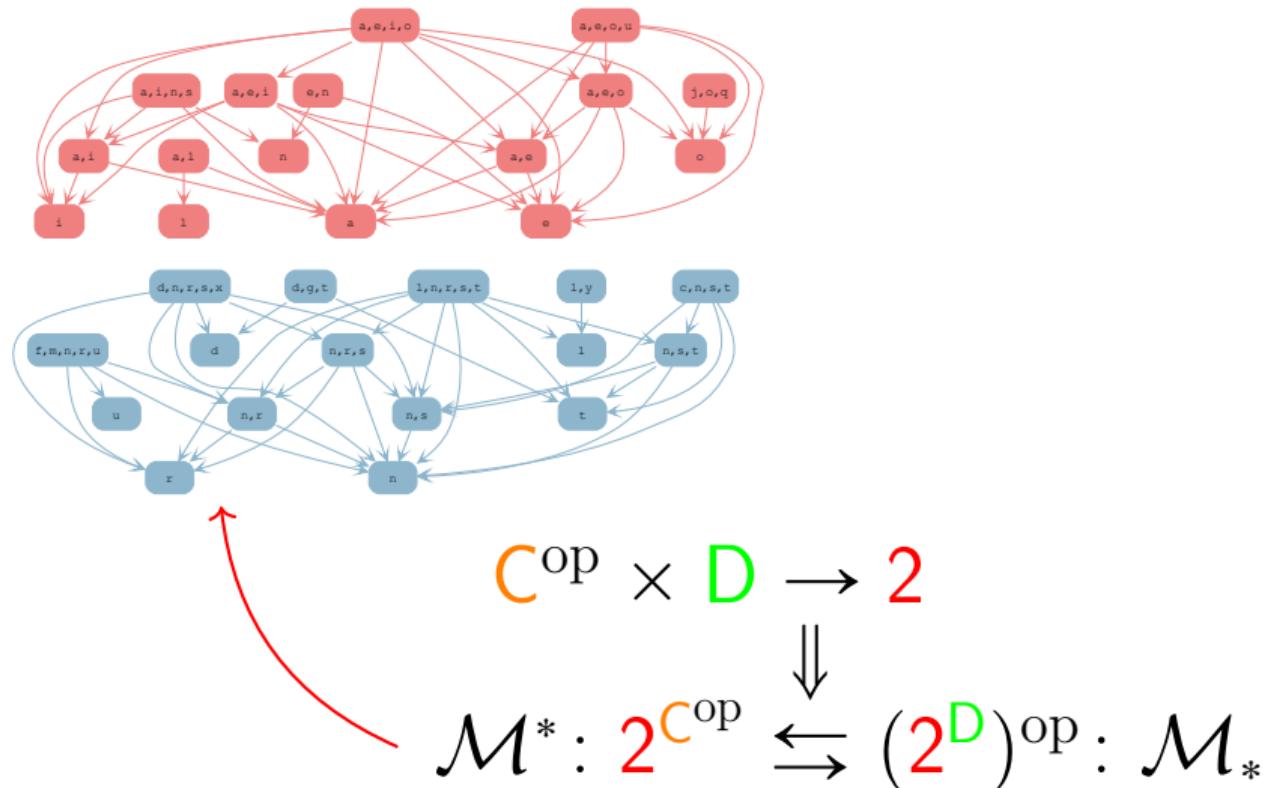


# Dual Partial Order



# Paring of Partial Ordered Fixed Points



Structure

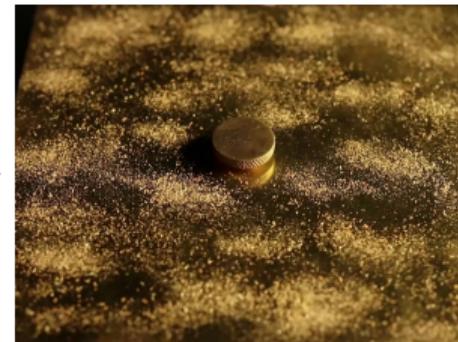
Structure

?

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

Structure

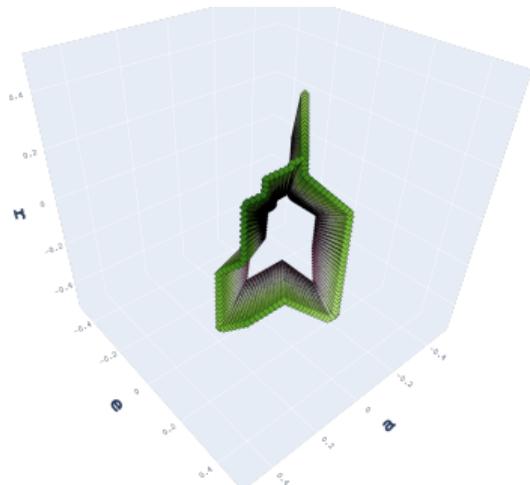
?



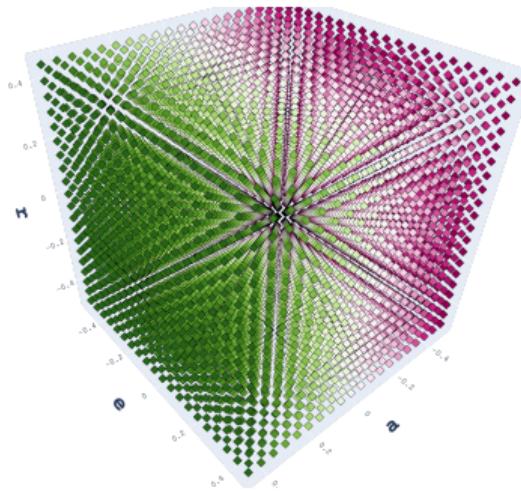
$$\begin{array}{c} \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}} \\ \Downarrow \\ \mathcal{M}^*: \bar{\mathbb{R}}^{\textcolor{orange}{C}^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^{\textcolor{green}{D}})^{\text{op}}: \mathcal{M}_* \end{array}$$

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

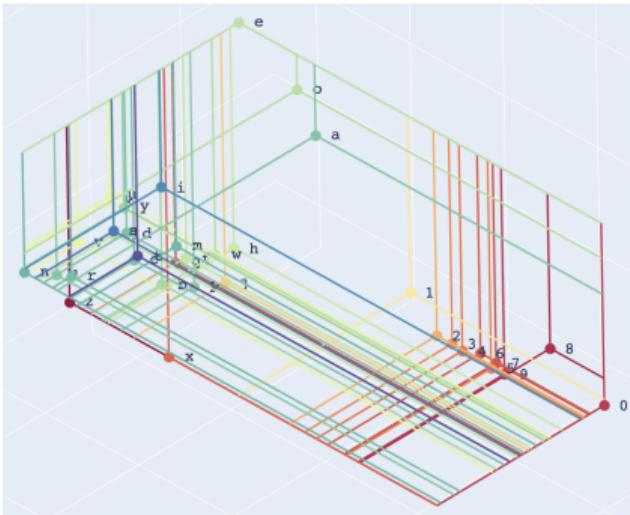
Structure



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



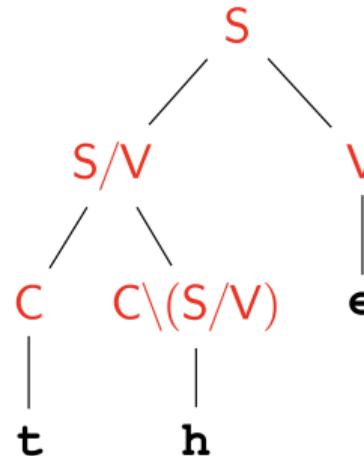
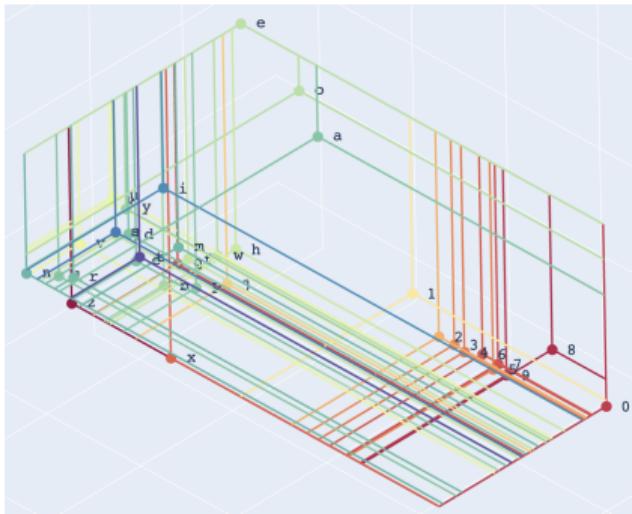
$$\begin{array}{c} \textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}} \\ \Downarrow \\ \mathcal{M}^*: \bar{\mathbb{R}}^{\textcolor{orange}{C}^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^{\textcolor{green}{D}})^{\text{op}}: \mathcal{M}_* \end{array}$$

Structure

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

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

## Structure

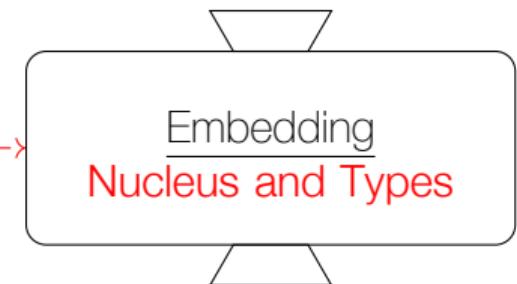
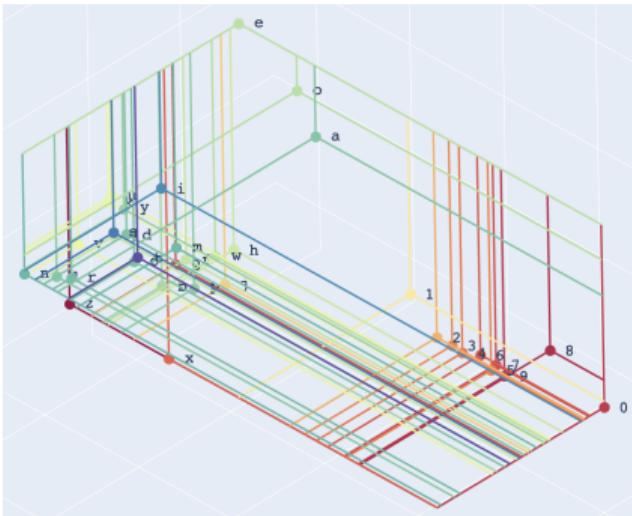


$$C^{\text{op}} \times D \rightarrow \bar{\mathbb{R}}$$



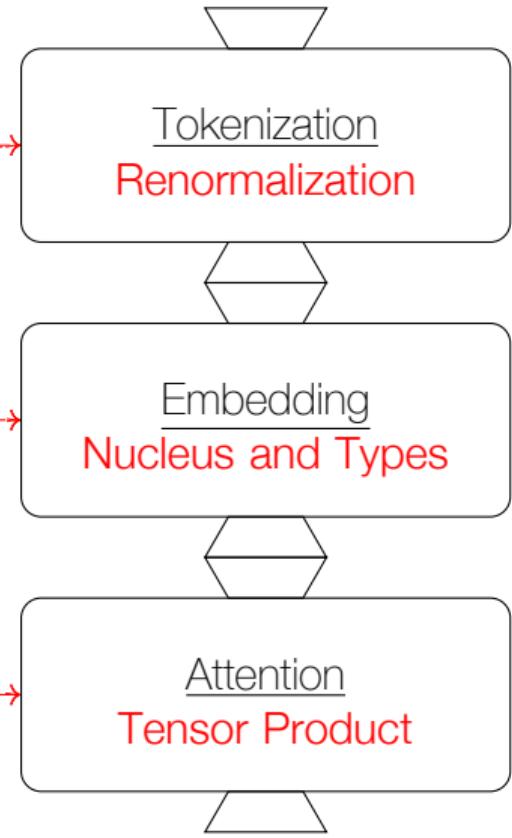
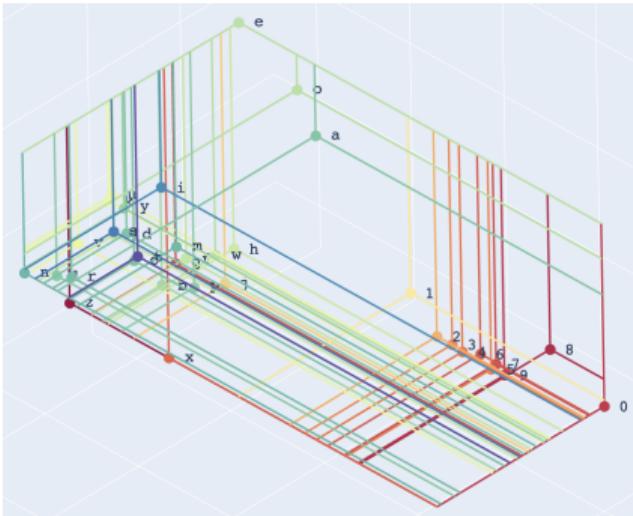
$$\mathcal{M}^*: \bar{\mathbb{R}}^{C^{\text{op}}} \rightleftarrows (\bar{\mathbb{R}}^D)^{\text{op}} : \mathcal{M}_*$$

## 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

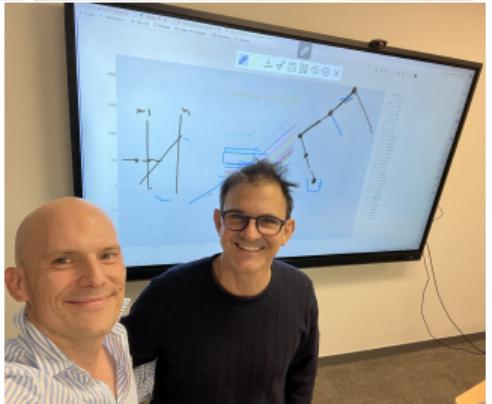
The Categories Behind the Structure

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
  - ◊ 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)

# Références |

- 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.
- 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., & 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.
- 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 catégories des cas*. Wilhelm Fink Verlag.

## Références II

- 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.
- 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.

## Références III

- 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/3f5ee243547dee91fb0d053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fb0d053c1c4a845aa-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

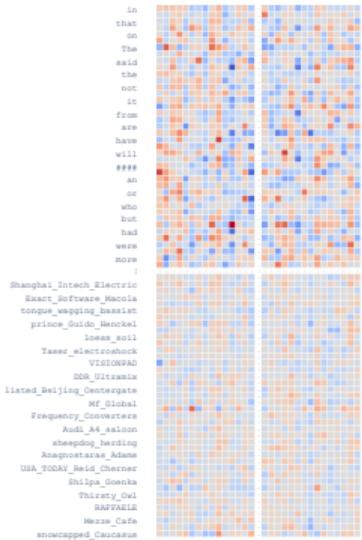
[www.giannigastaldi.com](http://www.giannigastaldi.com)

**ETH** zürich

June 12, 2025

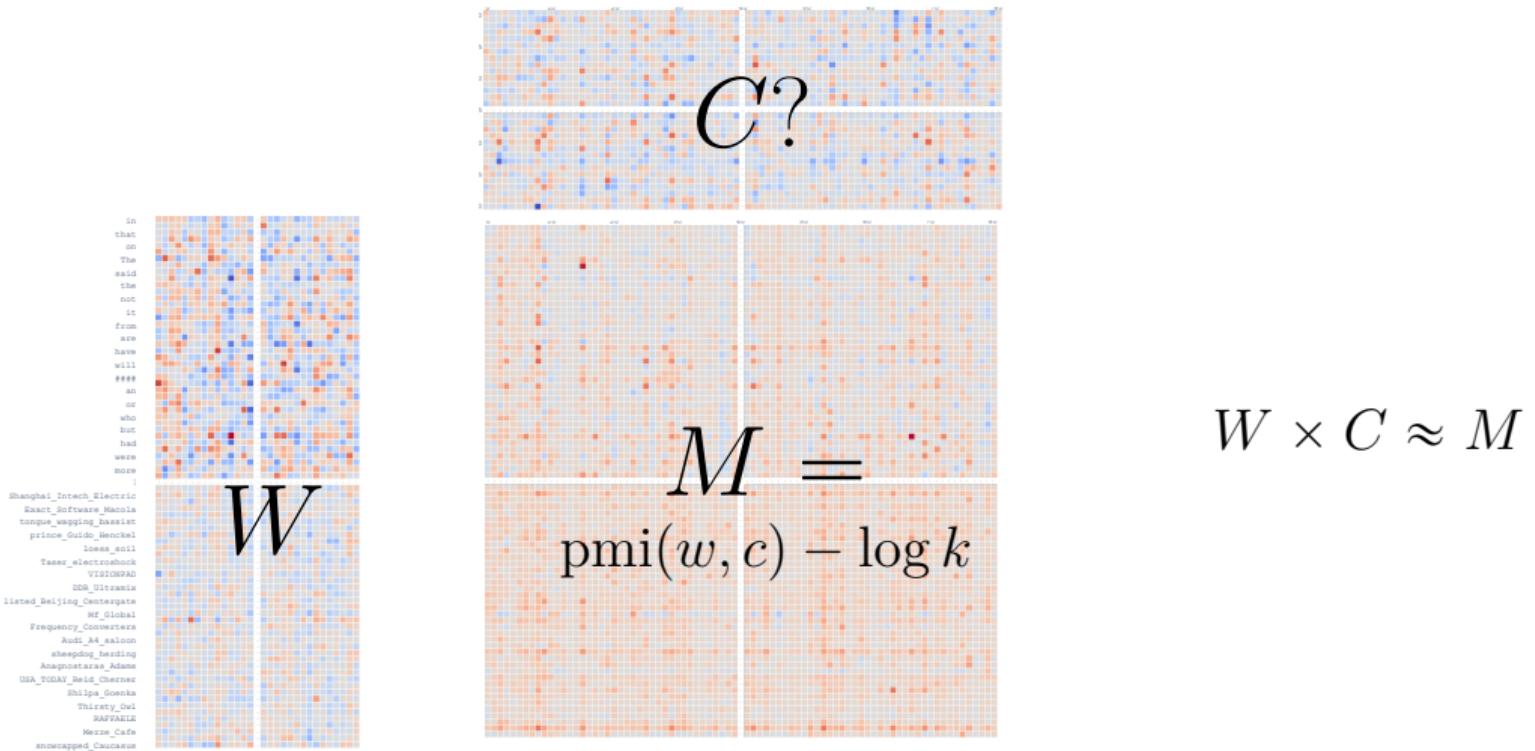
# word2vec as Implicit Matrix Factorization

(Levy and Goldberg, 2014)



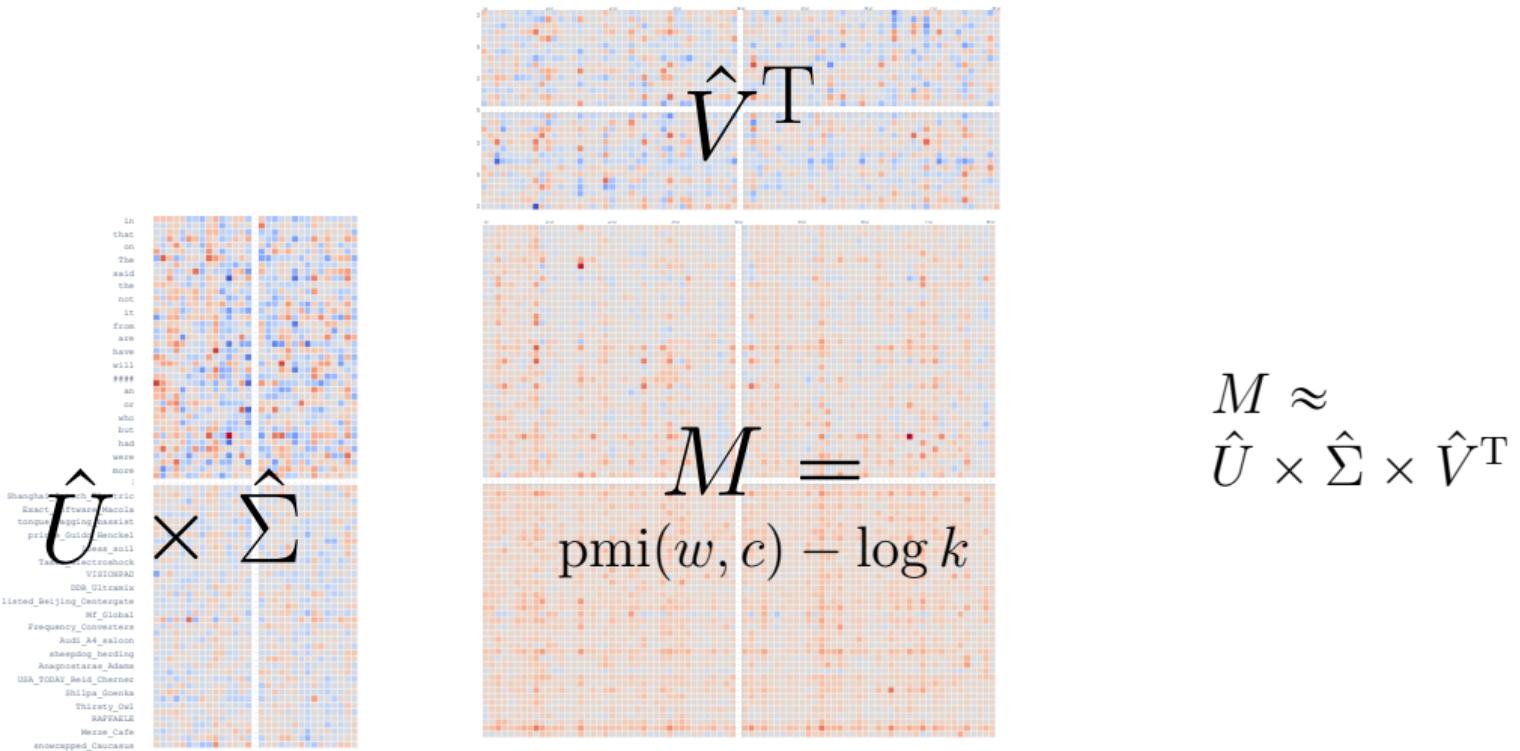
# word2vec as Implicit Matrix Factorization

(Levy and Goldberg, 2014)

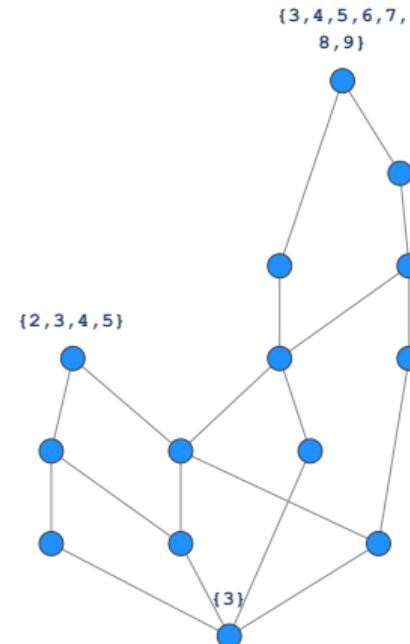
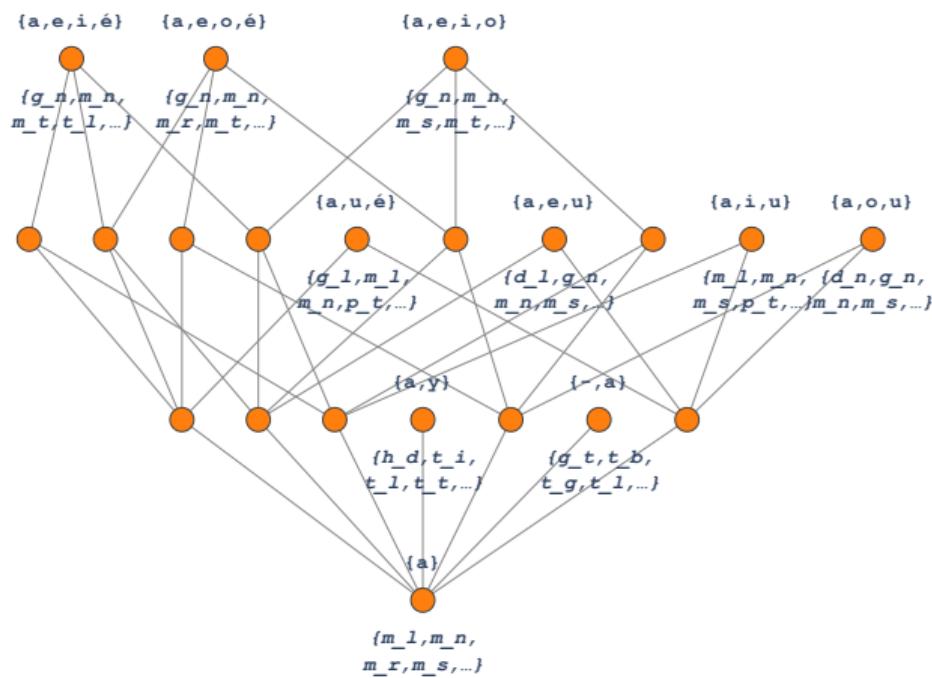


# Word Embeddings as Truncated SVD

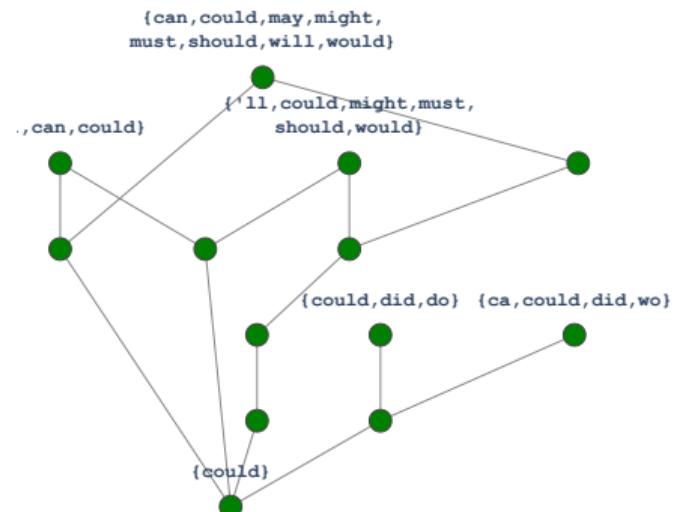
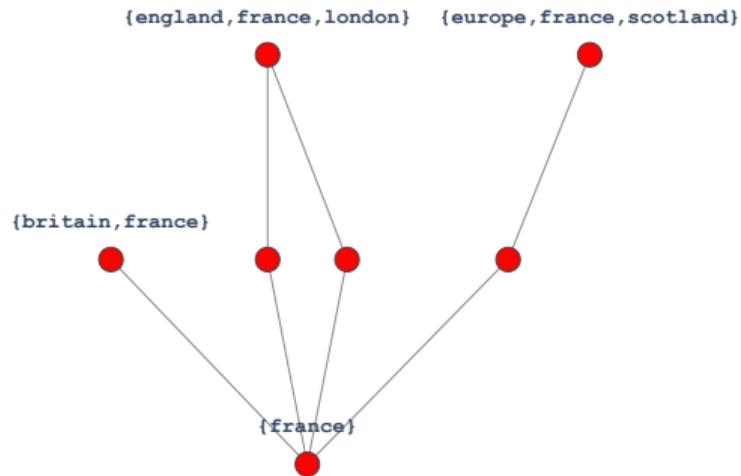
(Levy and Goldberg, 2014)



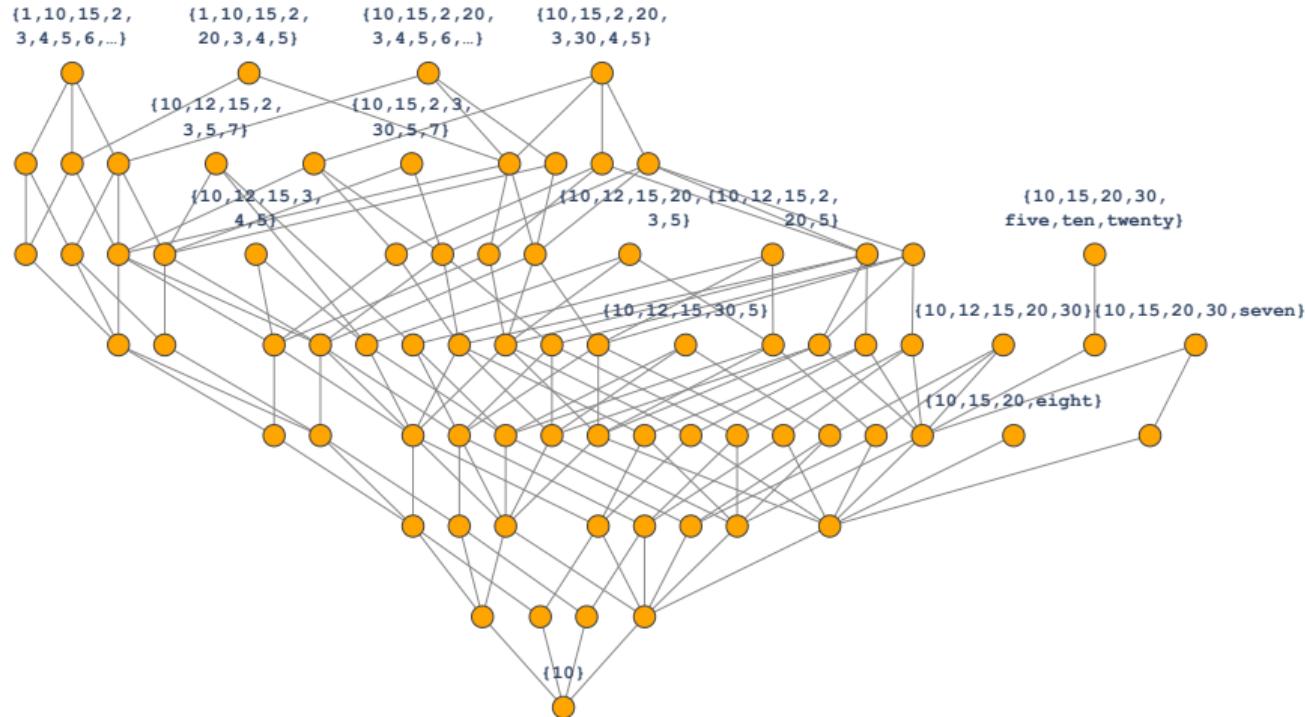
# Formal Concepts



# Formal Concepts (words)



# Formal Concepts (words)



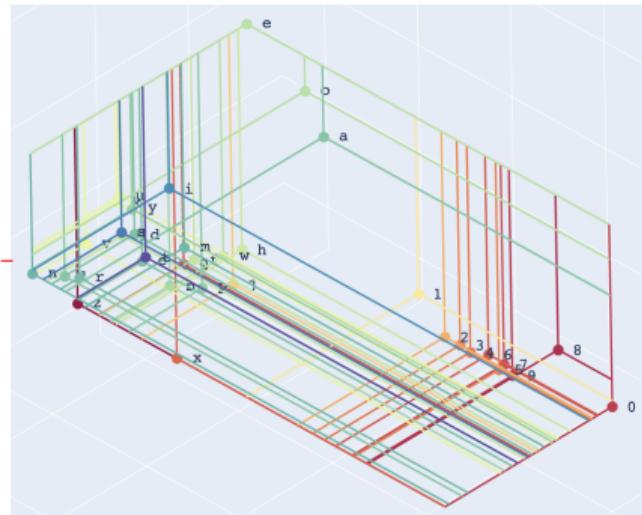
# Theoretical Interpretability

Theory  
"Task"

?



Structure



# Theoretical Interpretability

$$\textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}}$$

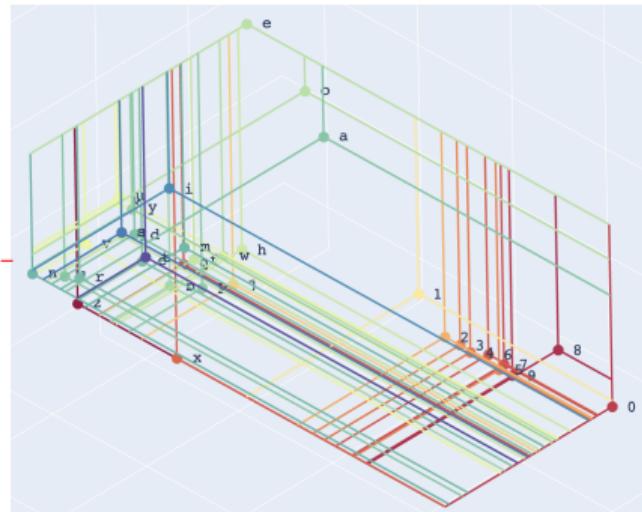
## Distributional Hypothesis

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

Theory  
“Task”



## Structure



# Theoretical Interpretability

$$\textcolor{orange}{C}^{\text{op}} \times \textcolor{green}{D} \rightarrow \bar{\mathbb{R}}$$

## Distributional Hypothesis

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

Theory  
“Task”

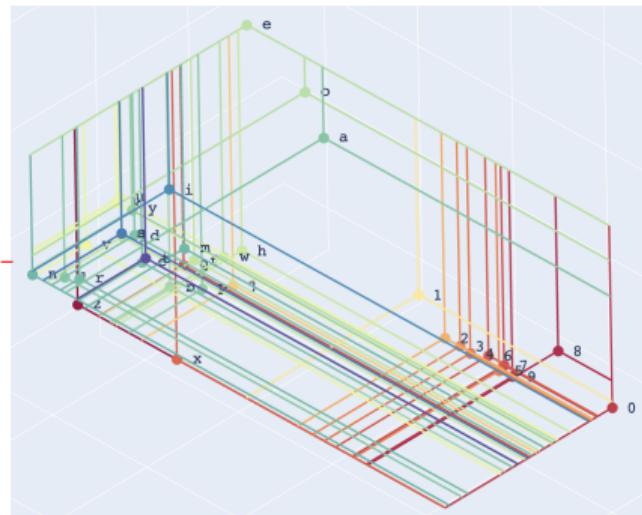


## Structuralist Hypothesis

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

$$\bar{\mathbb{R}}^{\textcolor{orange}{C}^{\text{op}}} \leftrightarrow (\bar{\mathbb{R}}^{\textcolor{green}{D}})^{\text{op}}$$

## Structure



# Theoretical Interpretability

Units

Classes

Relations

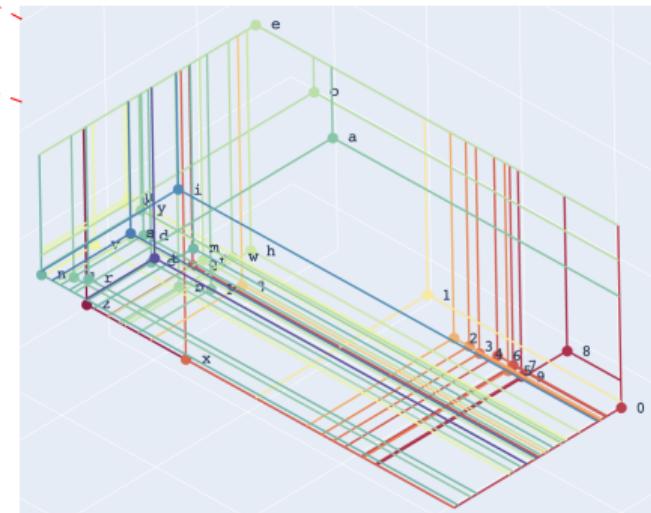
Semantics

Syntax

Morphology

Phonology

Structure



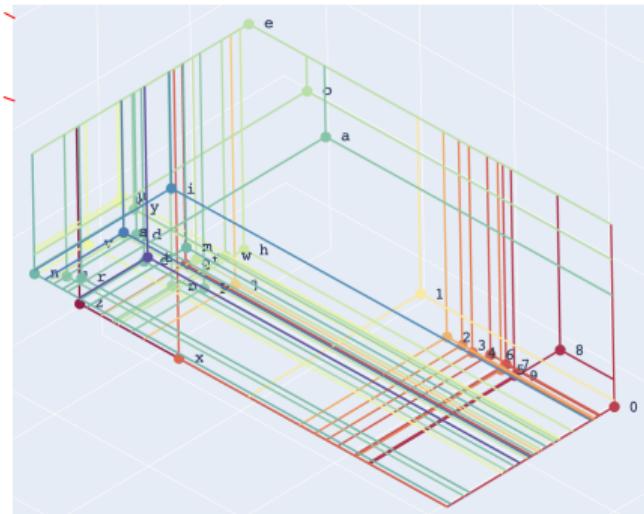
# Theoretical Interpretability

Units  
Classes  
Relations

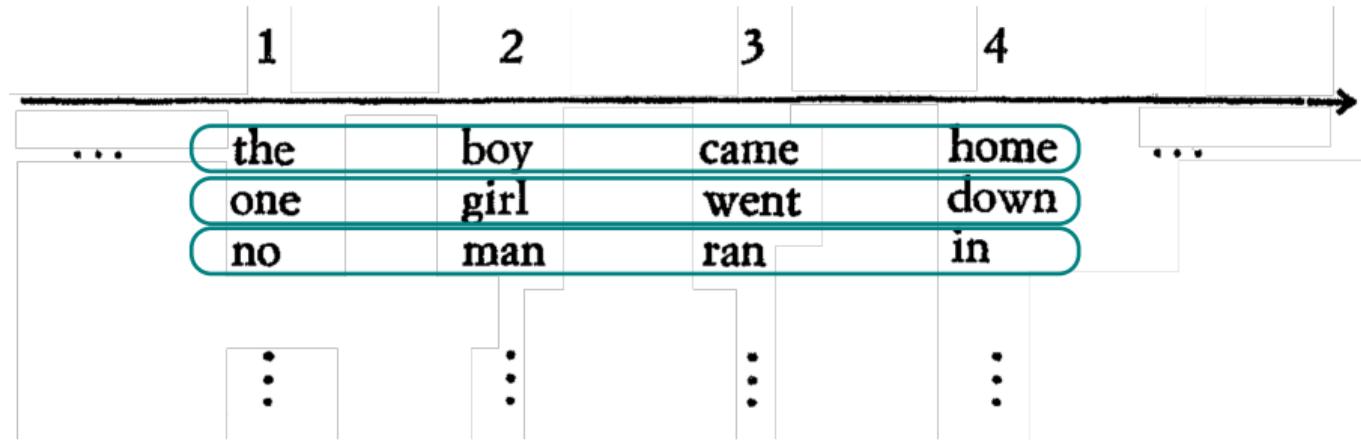
Semantics  
Syntax  
Morphology  
Phonology

	o	a	e	u	i	ɔ	ɛ	ɪ	ʊ	ə	ɔ̄	ɛ̄	ɪ̄	ʊ̄	f	ʃ	k	χ	g	ʒ	m	p	v	b	n	s	θ	t	z	ð	d	h	θ̄	h̄	#
1. Vocalic/Non-vocalic	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
2. Consonantal/Non-consonantal	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
3. Compact/Diffuse	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
4. Grave/Acute	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
5. Flat/Plain	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
6. Nasal/Oral	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
7. Tense/Lax	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
8. Continuant/Interrupted	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
9. Strident/Mellow	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		

(Jakobson et al., 1952)

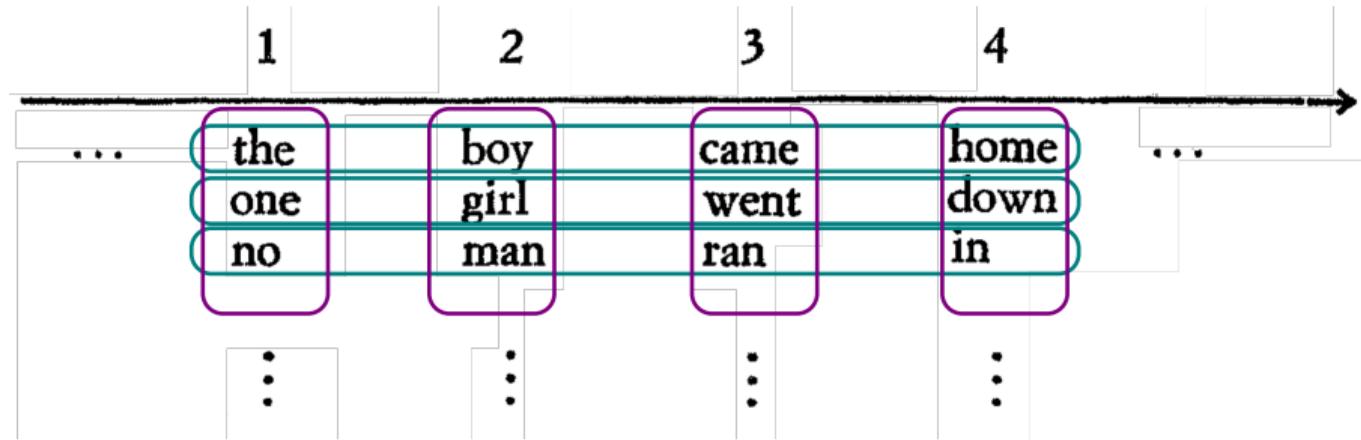


# Syntagmes et Paradigmes



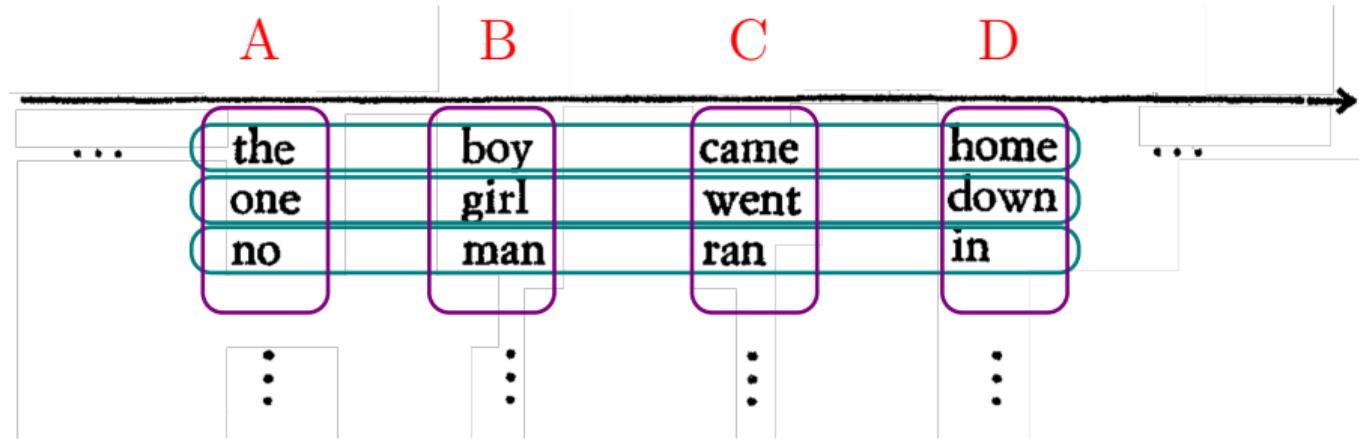
(Hjelmslev, 1971)

# Syntagmes et Paradigmes



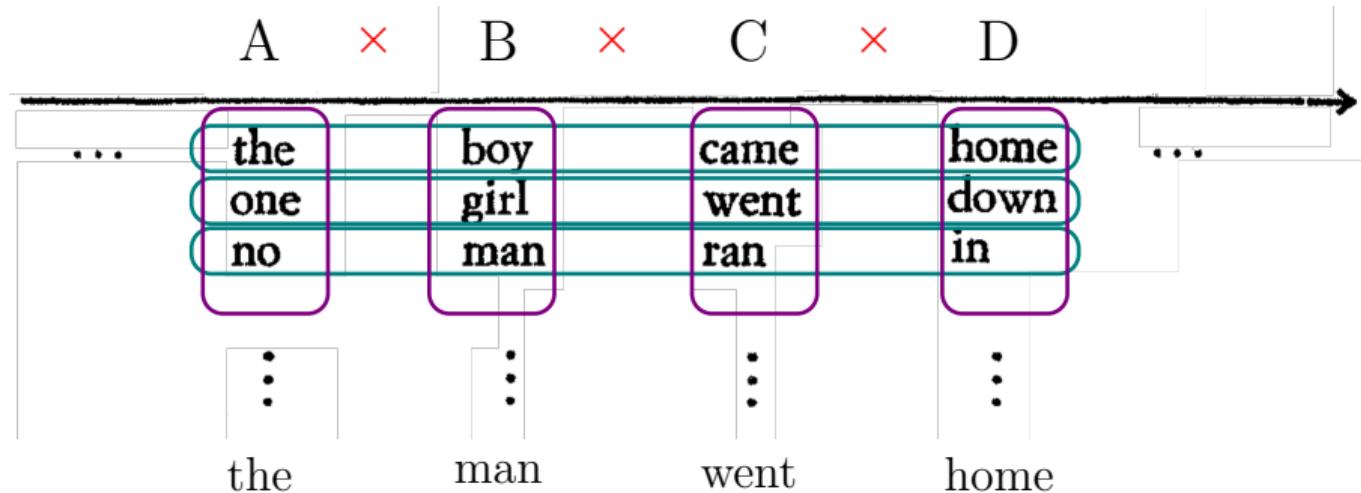
(Hjelmslev, 1971)

# Syntagmes et Paradigmes



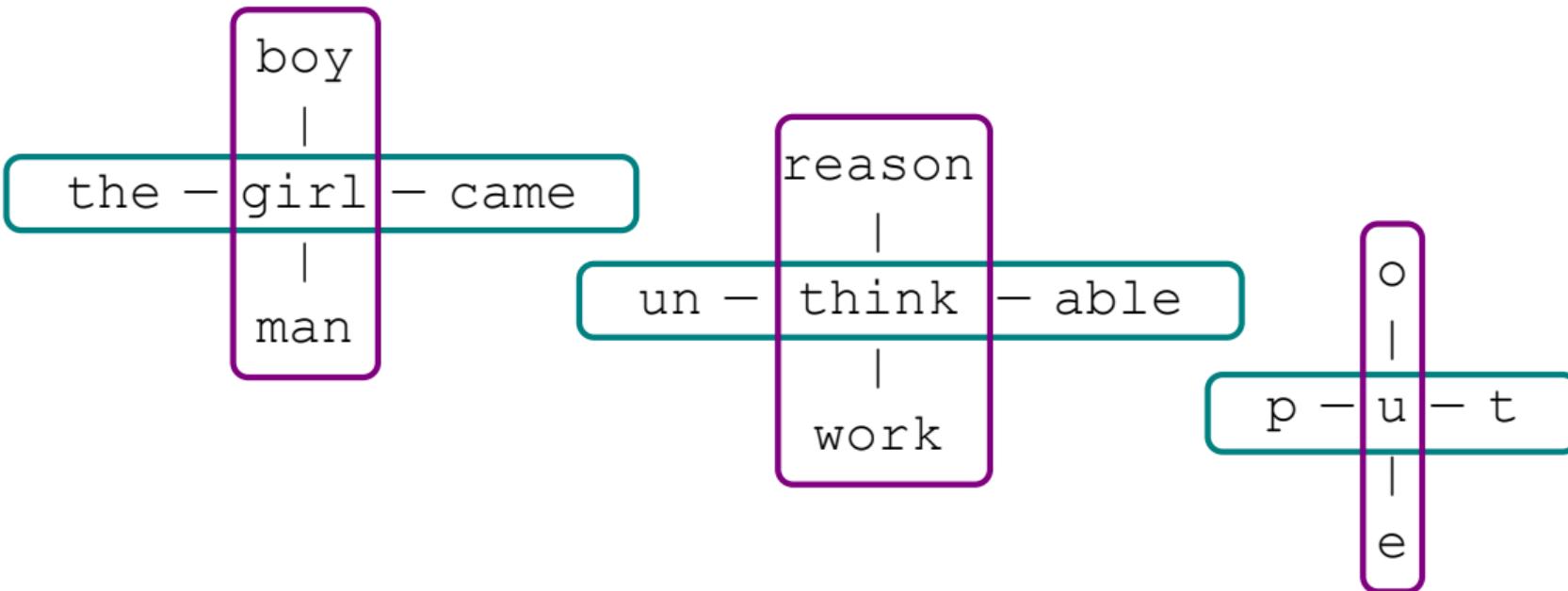
(Hjelmslev, 1971)

# Syntagmes et Paradigmes

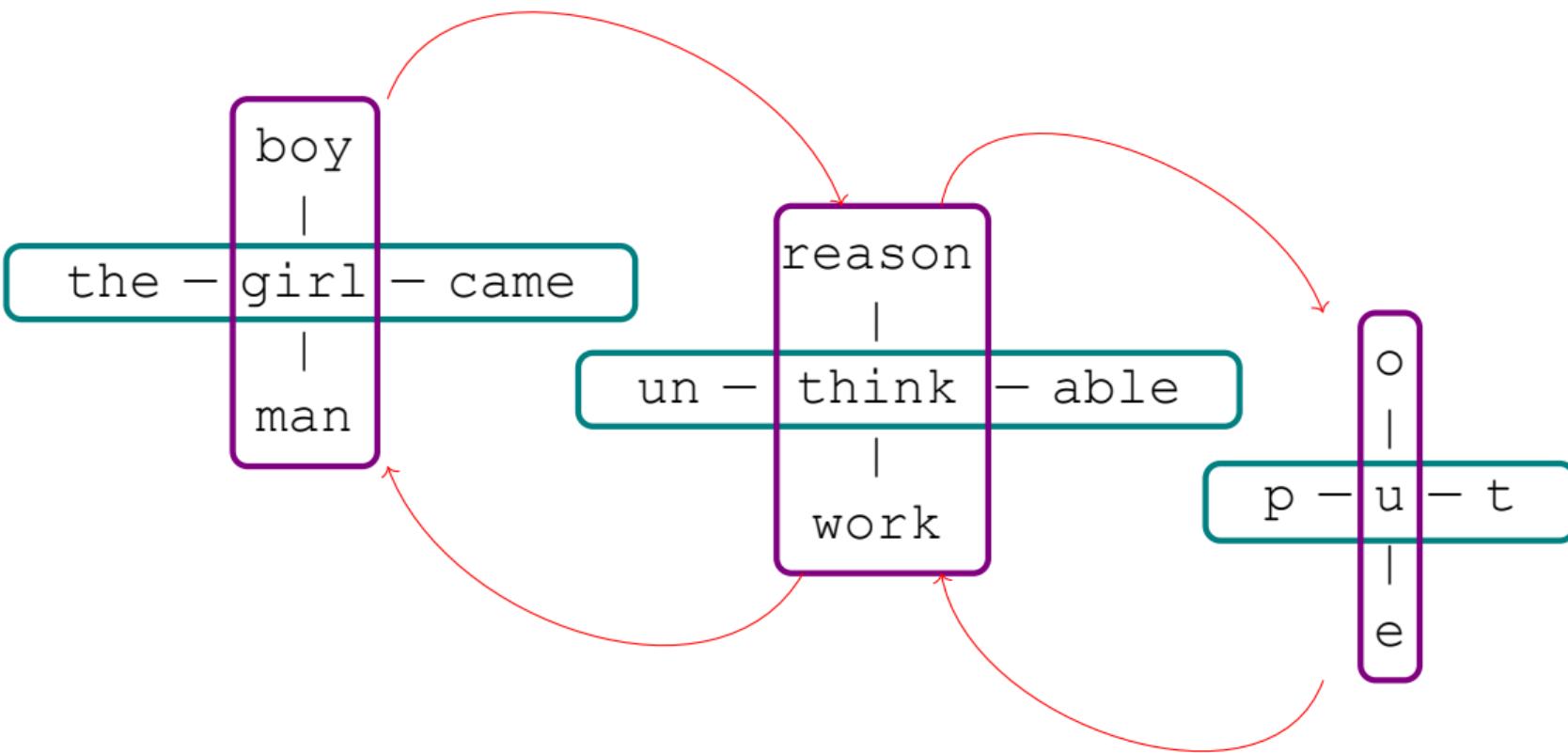


(Hjelmslev, 1971)

# Stratification



# Stratification



# De l'algèbre linéaire aux catégories

$$\begin{array}{ccc}
 X & \xrightarrow{M_x} & \mathbb{R}^Y \\
 \downarrow & \nearrow M^* & \uparrow \\
 \mathbb{R}^X & \xleftarrow{M_y} & Y
 \end{array}$$

$$\begin{array}{ccc}
 C & \xrightarrow{\mathcal{M}_c} & (\text{Set}^D)^{\text{op}} \\
 \downarrow \text{Yoneda} & \nearrow \mathcal{M}^* & \uparrow \text{Yoneda} \\
 \text{Set}^{C^{\text{op}}} & \xleftarrow{\mathcal{M}_d} & D
 \end{array}$$

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

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

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

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

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

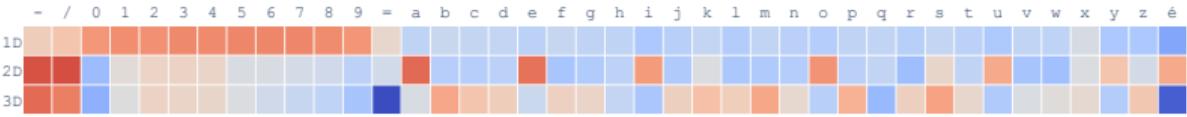
$$\mathcal{M}^* \mathcal{M}_*: (\text{Set}^D)^{\text{op}} \rightarrow (\text{Set}^D)^{\text{op}}$$

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

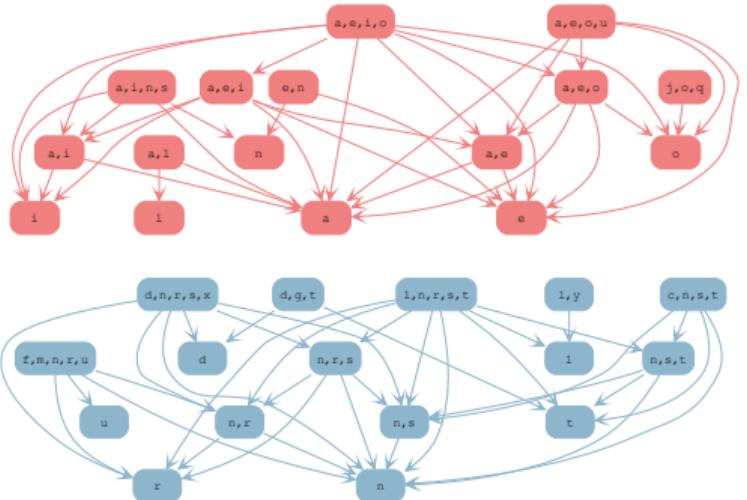
$$\text{Fix}(\mathcal{M}^* \mathcal{M}_*) := \{g \in (\text{Set}^D)^{\text{op}} \mid \mathcal{M}^* \mathcal{M}_*(g) \cong g\}$$

# Structures catégoriques

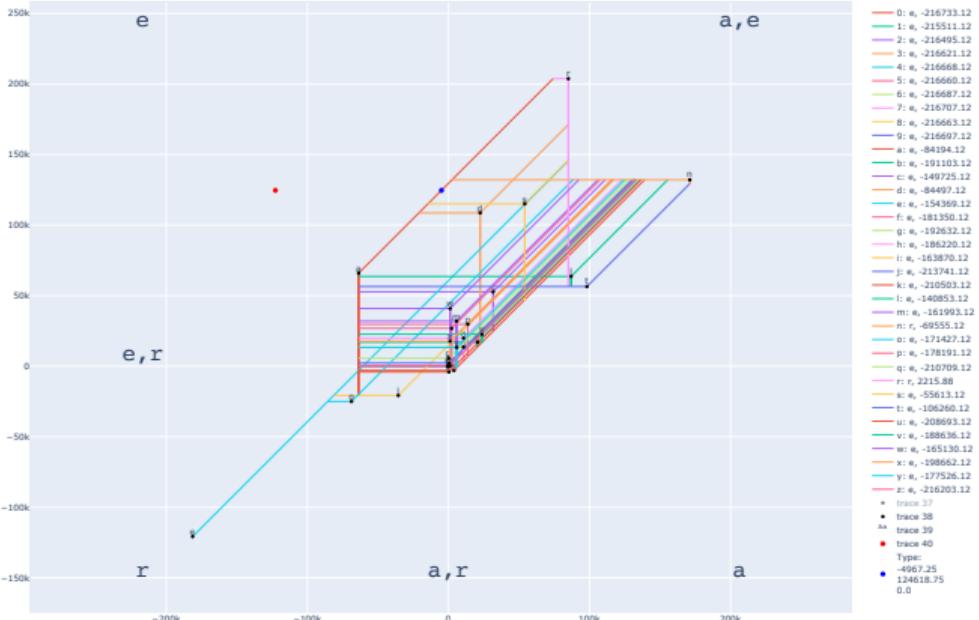
## Algèbre linéaire



$$\mathcal{M}^*: \textcolor{red}{2^C}^{\text{op}} \leftrightarrows (\textcolor{red}{2^D})^{\text{op}}: \mathcal{M}_*$$



$$\mathcal{M}^*: \bar{\mathbb{R}}^{\text{C}^{\text{op}}} \leftrightarrows (\bar{\mathbb{R}}^{\text{D}})^{\text{op}}: \mathcal{M}_*$$



# Théorie des types computationnels

## Definition (Polaire/Orthogonal - Girard, 2006)

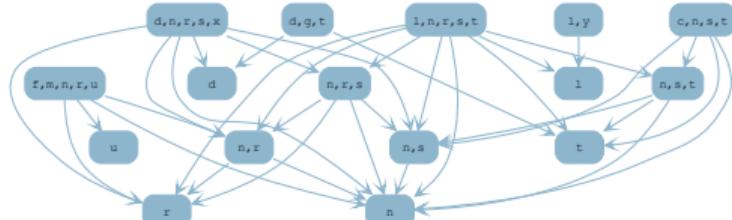
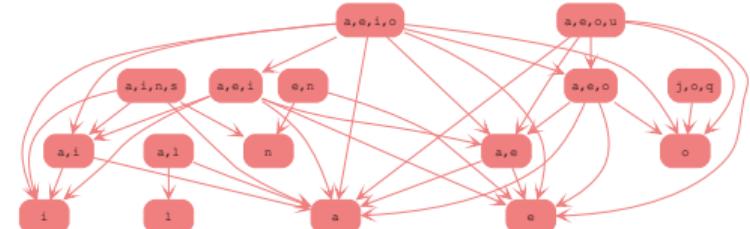
[É]tant donnée une fonction binaire

$a, b \rightsquigarrow \langle a|b \rangle : A \times B \rightarrow C$  et un sous-ensemble  $P \subset C$  (le « pôle »), on peut définir le *polaire*  $X^\perp \subset B$  d'un sous-ensemble  $X \subset A$  (resp.  $Y^\perp \subset A$  d'un sous-ensemble  $Y \subset B$ ) par :

$$X^\perp := \{y \in B : \forall x \in X, \langle a|b \rangle \in P\}$$

$$Y^\perp := \{x \in A : \forall y \in Y, \langle a|b \rangle \in P\}$$

- ◊ L'application « polaire » est décroissante:  $X \subset X' \Rightarrow X'^\perp \subset X^\perp$ .
- ◊ L'ensemble  $\text{Pol}(A) \subset \mathcal{P}(A)$  des ensembles *polaires*, i.e., de la forme  $Y^\perp$ , est stable par intersections arbitraires. En particulier,  $A$  est polaire et  $X^{\perp\perp}$  est le plus petit ensemble polaire contenant  $X$ .
- ◊ En conséquence,  $X^{\perp\perp\perp} = X^\perp$ .



# Axe 2: Interprétabilité théorique

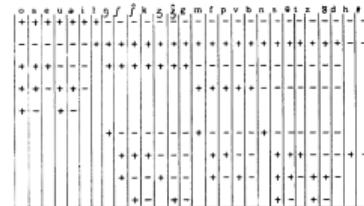
## Hypothèse distributionnelle

Le contenu des unités linguistiques est déterminé par leur distribution dans un corpus.

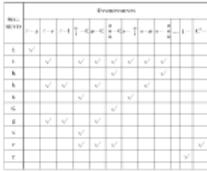


(Hjelmslev, 1935)

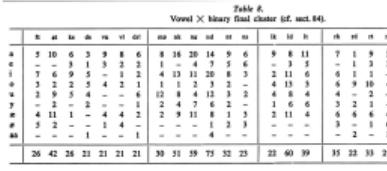
1. Vocalic/Non-vocalic
2. Consonantal/Non-consonantal
3. Compact/Diffuse
4. Grave/Acute
5. Flat/Plain
6. Nasal/Oral
7. Tense/Lax
8. Continuant/interrupted
9. Strident/Mellow



(Jakobson et al., 1952)



(Harris, 1960)



(Spang-Hanssen, 1959)

## Hypothèse structurale

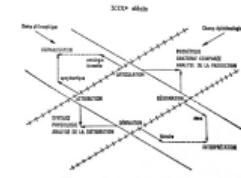
Le contenu linguistique est l'effet d'une structure virtuelle dérivée des pratiques linguistiques dans une communauté.

Répondeur	Fonc. fonctionnel	fonc. fonctionnel	fonc. fonctionnel
Prix	+	+	-
Prix, idéal	-	-	-
Biens	+	+	+
Prise de la parole	+	+	+
Signification de l'acte	-	-	-
Acte	+	+	+
Cod	+	+	+
Progrès	+	+	+
Signification de l'acte	-	-	-
Biens à des fins extrinsèques	+	+	+
Autonomie	+	+	+
Obtention	+	+	+
Signification de l'acte	-	-	-
Progrès de l'acte	+	+	+
Signification de l'acte	-	-	-

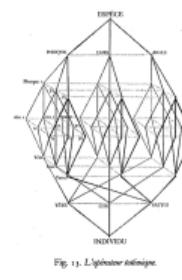
(Lévi-Strauss, 1949)



(Bourdieu, 1979)



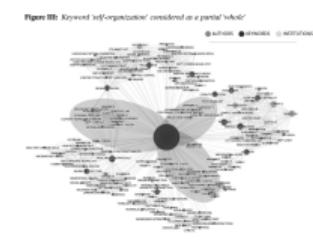
(Foucault, 1966)



(Lévi-Strauss, 1962)

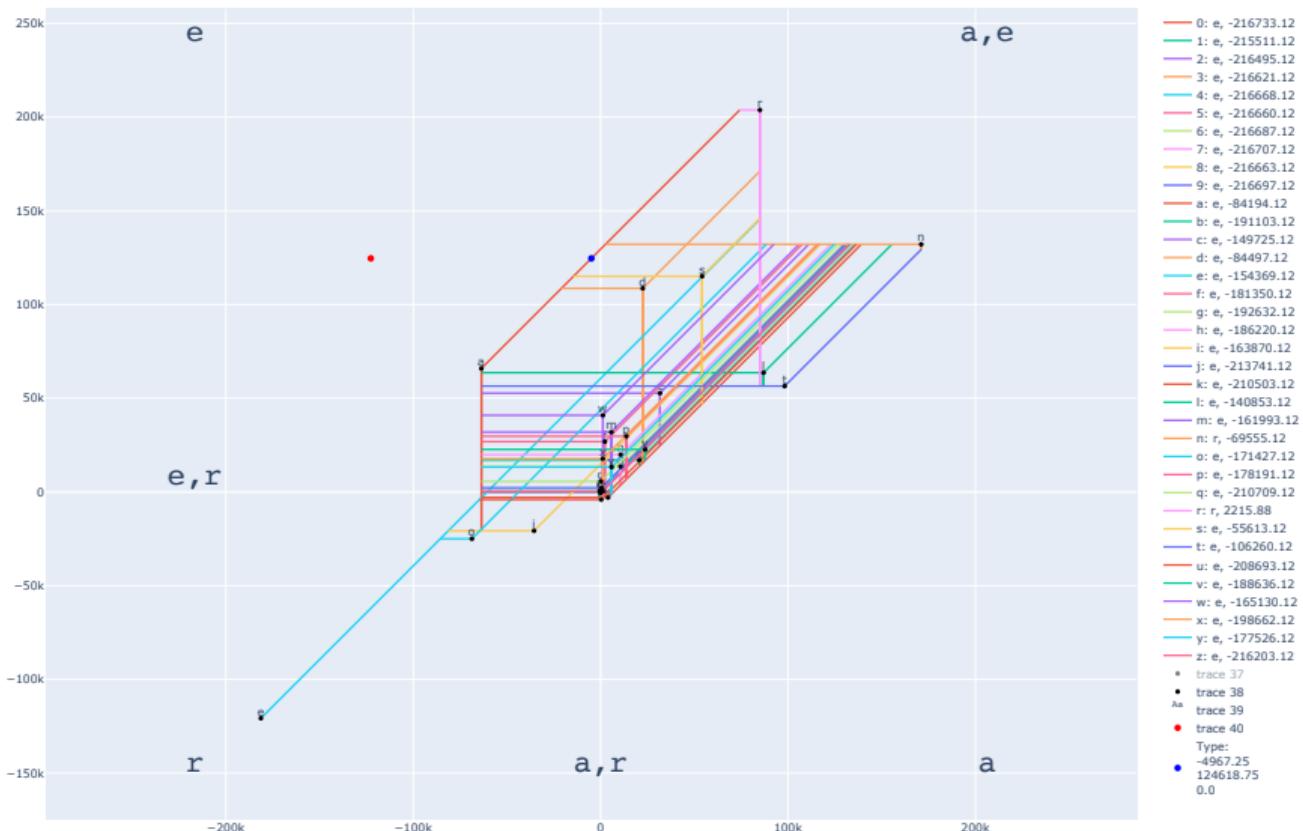


(Bourdieu, 1994)

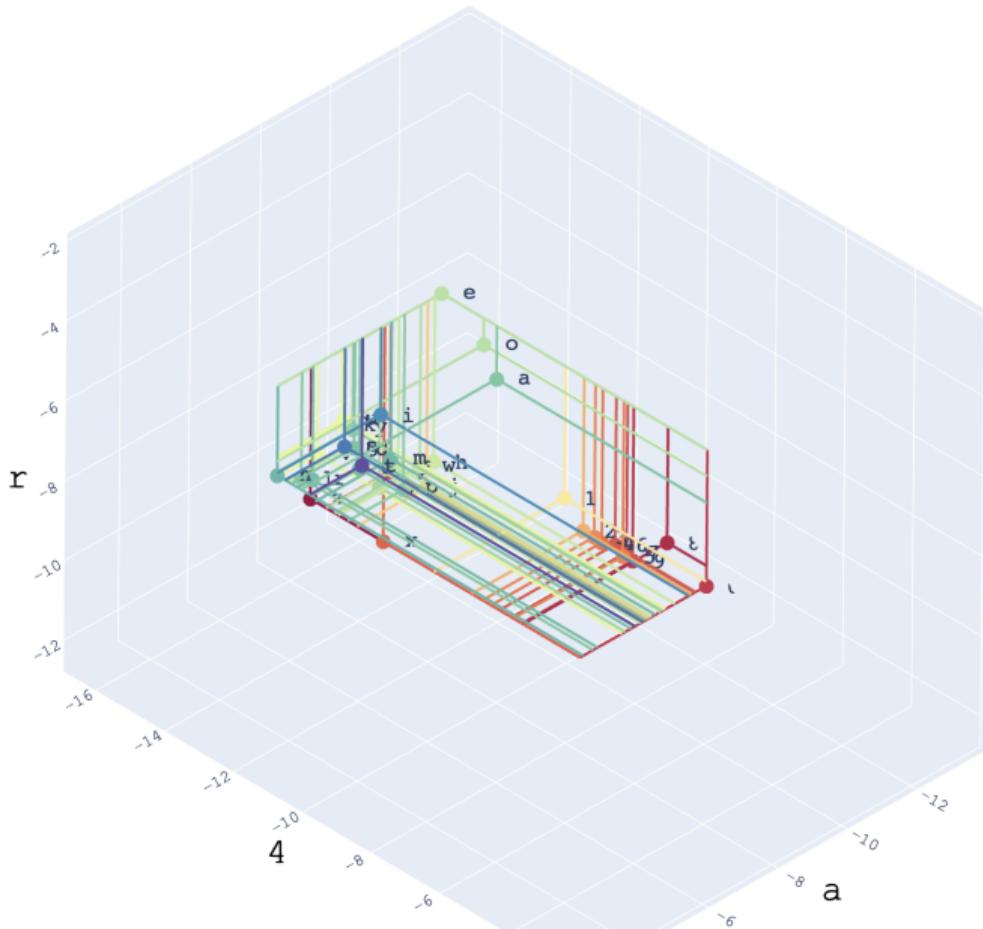


(Latour et al., 2012)

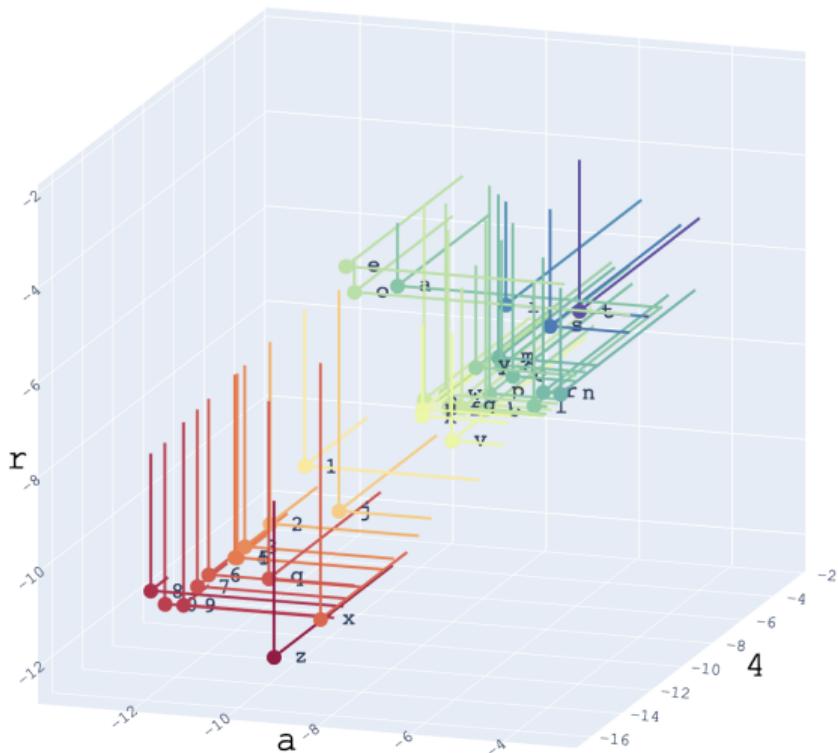
# Structure interne du noyau



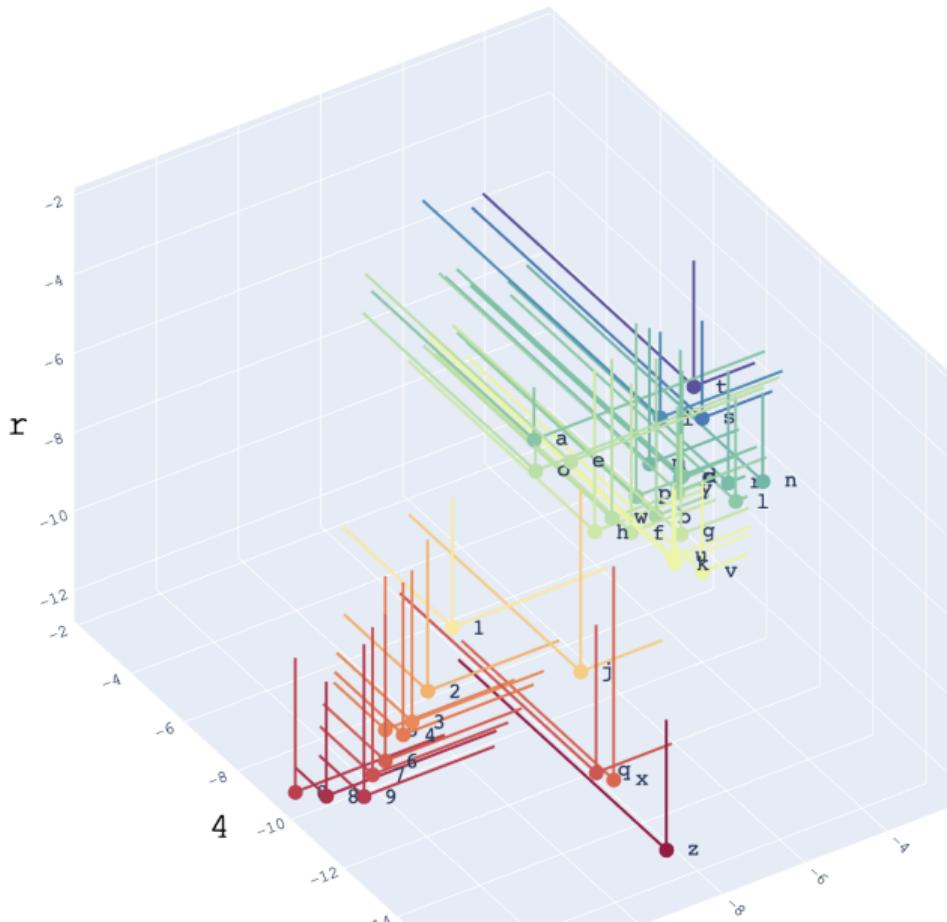
# Structure interne du noyau



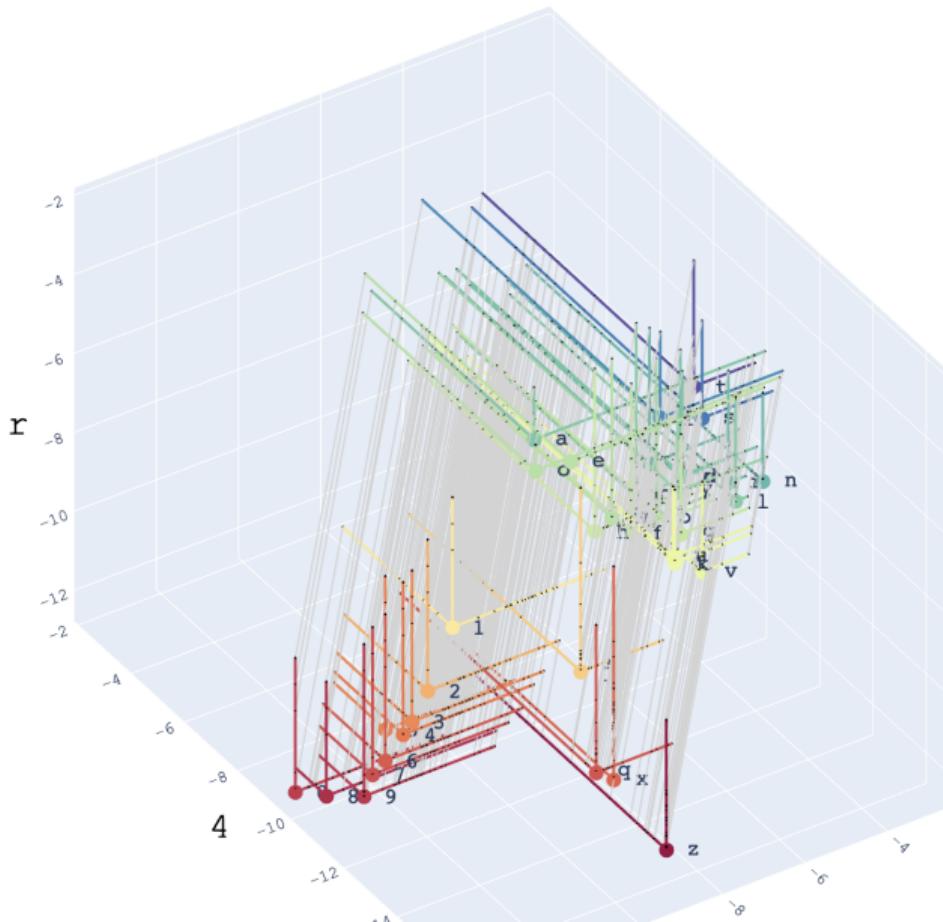
# Structure interne du noyau



# Structure interne du noyau



# Structure interne du noyau



# Matrice et analogie

a = your  
c = my

w = apartment  
x = house  
y = chair  
z = stool

your : house  
my : apartment

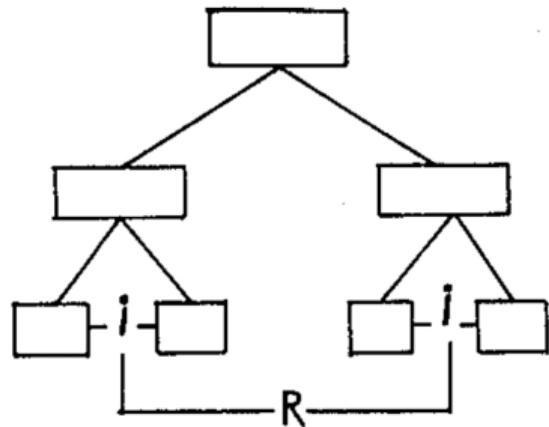
	...	w	x	y	z	...
...	...	0	0	0	0	...
a	...	0	1	1	0	...
b	...	0	0	1	1	...
c	...	1	0	0	1	...
...	...	0	0	0	0	...

Une **sémiotique** [...] est une hiérarchie dont chacune des composantes admet une analyse ultérieure en classes définies par relation mutuelle, de telle sorte que chacune de ces classes admette une analyse en dérivés définis par mutation mutuelle.

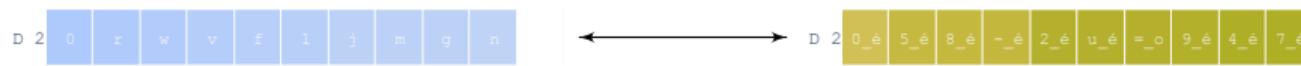
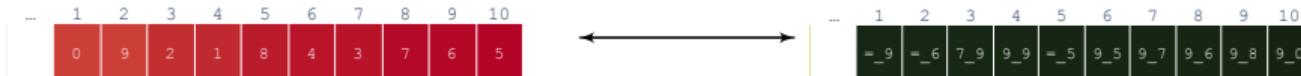
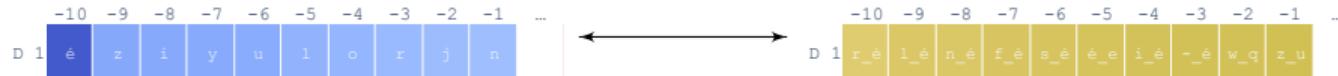
(Hjelmslev, 1975, Df. 24)

Une **mutation** [...] est une fonction existant entre des dérivés du premier degré d'une seule et même classe, une fonction qui a une relation à une fonction entre d'autres dérivés de premier degré d'une seule et même classe et appartenant au même rang.

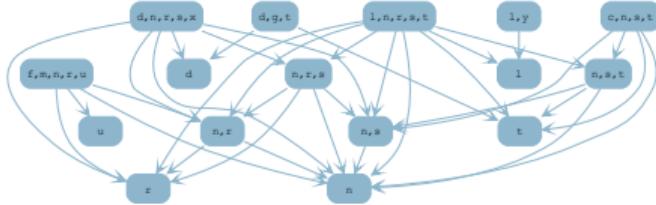
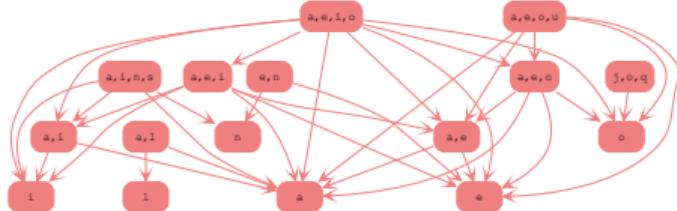
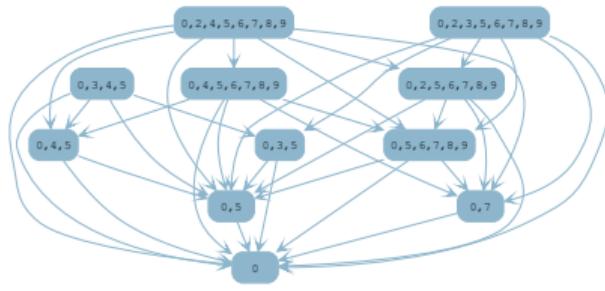
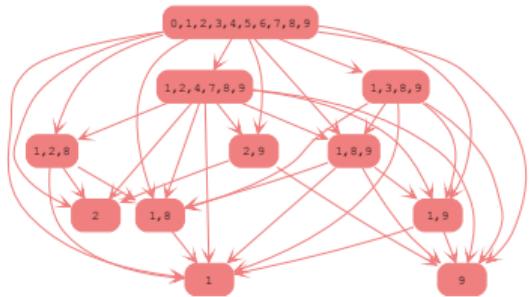
(Hjelmslev, 1975, Df. 23)



## Syntagmatique et Texte (Vecteurs)



# Syntagmatique et Texte (Noyau/Types)



# Paradigmatique et Langue (Vecteurs)

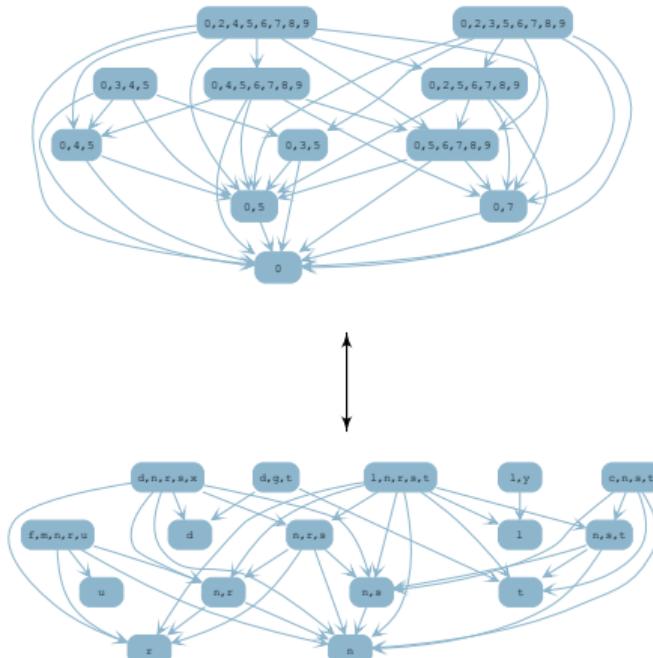
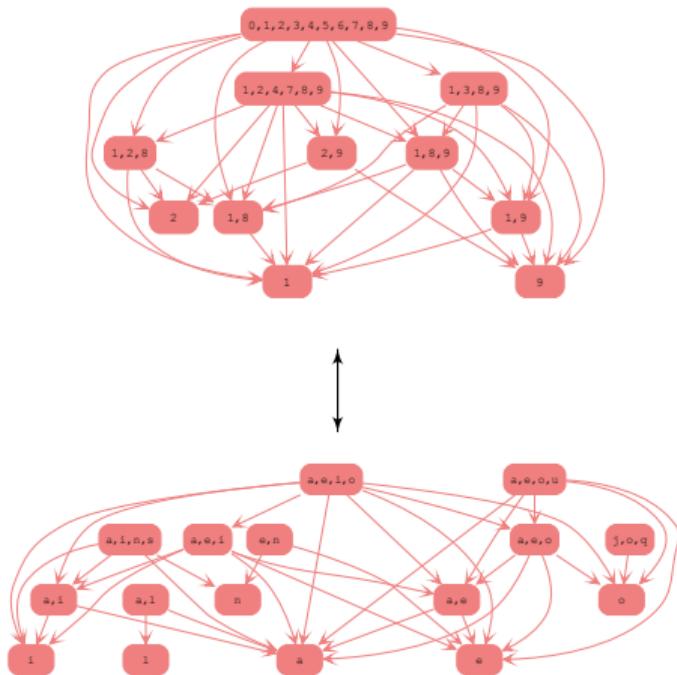
-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	...
D 1	é	z	i	y	u	l	o	r	j	n
...	1	2	3	4	5	6	7	8	9	10
	0	9	2	1	8	4	3	7	6	5

-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	...
D 1	x_é	l_é	n_é	f_é	s_é	e_é	i_é	u_é	w_q	z_u
...	1	2	3	4	5	6	7	8	9	10
	=_9	=_6	7_-9	9_-9	=_5	9_-5	9_-7	9_-6	9_-8	9_0

0	r	w	v	f	l	j	m	g	n	...
D 2	3	y	u	é	i	o	e	a	-	/
	z	p	f	m	g	t_g	é_m	z_m	z_g	z_q

0_é	5_é	8_é	-_é	2_é	u_é	=_o	9_é	4_é	7_é	...
D 2	d_m	z_p	z_f	k_m	r_g	t_g	é_m	z_m	z_g	z_q
	z	p	f	m	g	t_g	é_m	z_m	z_g	z_q

# Paradigmatique et Langue (Noyau/Types)



# Illustration du contenu formel

(Gastaldi and Pellissier, 2021)

## Characteristic Content

```
{cat, dog, spider,  
gavagai}
```

Atomic Type

## Syntactic Content

"the gavagai is on the  
mat"

Profunctor Nucleus

## Informational Content

```
{cat:0.059%,  
dog:0.012%,  
spider:0.009%  
gavagai:0.000%}
```

Probability Distribution