

Vernacular AI
Digital Theory Lab, NYU
New York, NY, USA

The Structure, Not the Prompt
For a Critical Formalism

Juan Luis Gastaldi

ETH zürich

February 7, 2025

Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

 The Algebra Behind the Embeddings

 The Structure Behind the Algebra

 The Categories Behind the Structure

Conclusion

Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

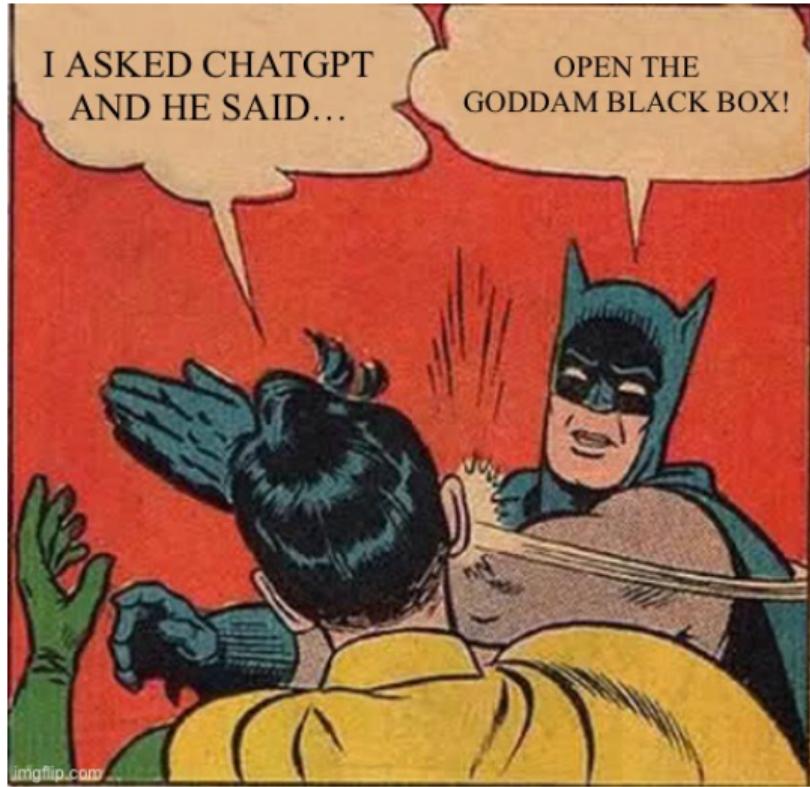
The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

...Not the Prompt



Where Art Thou, 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)
 - 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 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)

Critique and Formalism

- ◊ 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...

Critique and Formalism

- ◊ 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 and Formalism

- ◊ 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**

Critique and Formalism

- ◊ 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

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Neural LMs as Computable Functions

Neural LM



?

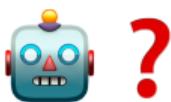
Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions

Neural LM

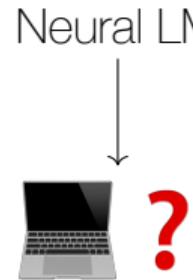


Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions



Neural LMs as Computable Functions

Neural LM



Neural LMs as Computable Functions

Neural LM



$$f !$$

Neural LMs as Computable Functions

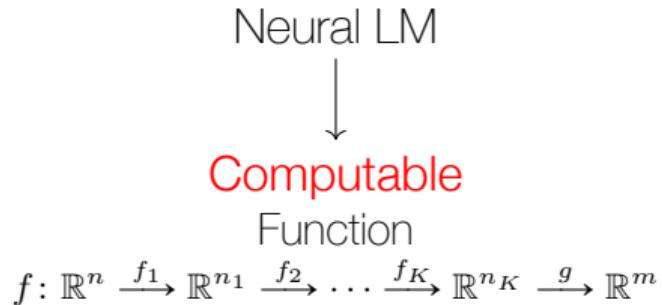
Neural LM



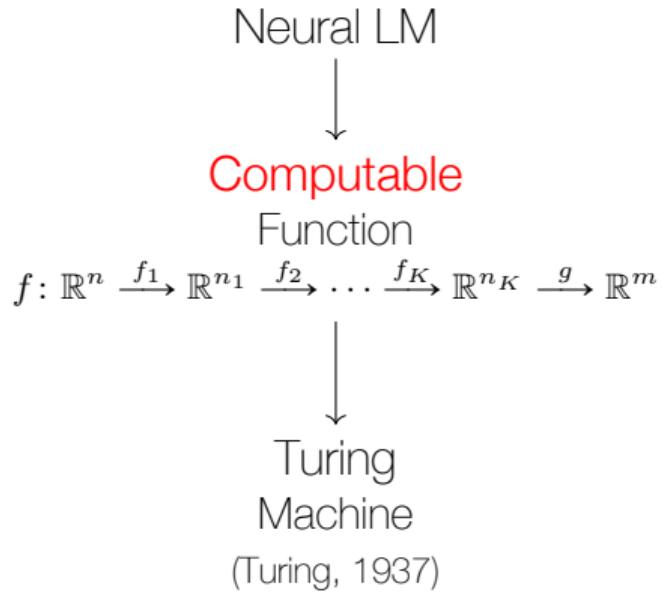
Function

$$f: \mathbb{R}^n \xrightarrow{f_1} \mathbb{R}^{n_1} \xrightarrow{f_2} \dots \xrightarrow{f_K} \mathbb{R}^{n_K} \xrightarrow{g} \mathbb{R}^m$$

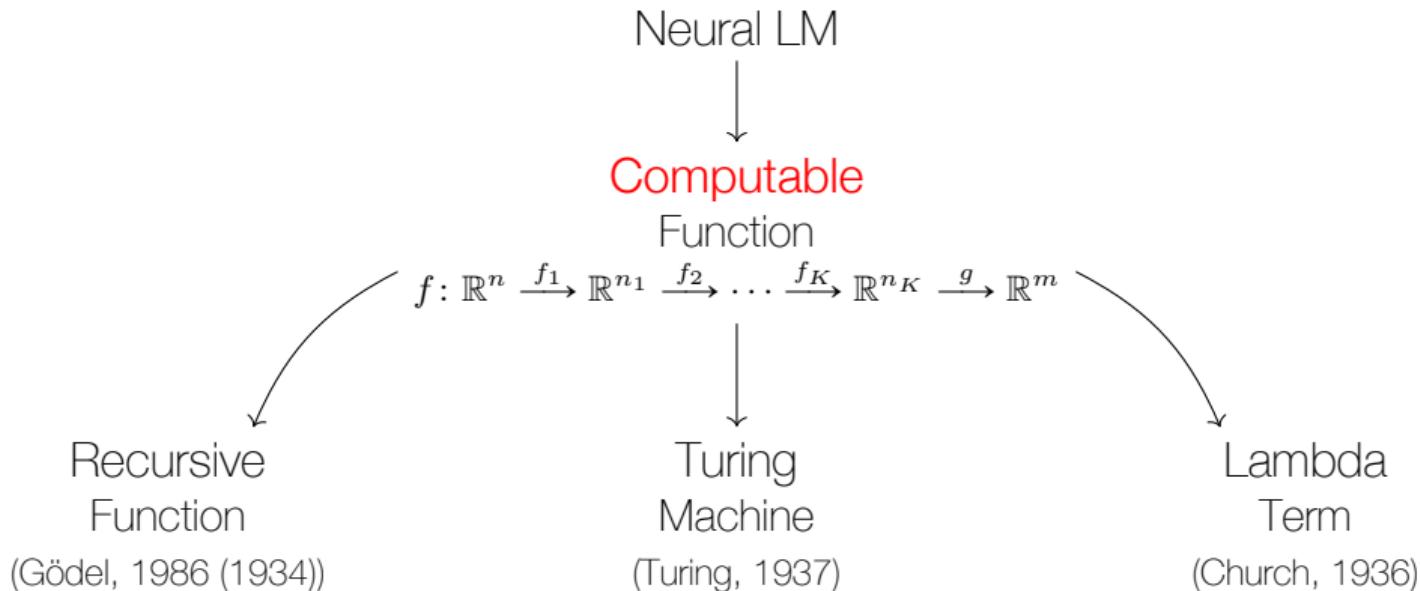
Neural LMs as Computable Functions



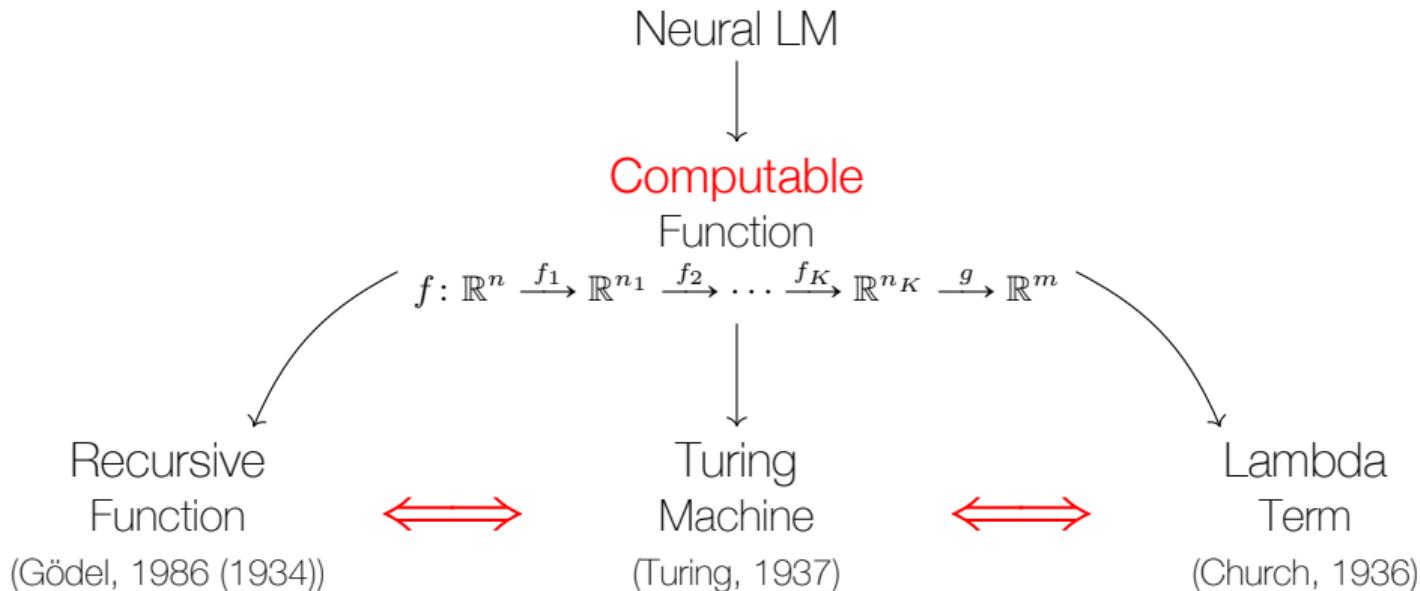
Neural LMs as Computable Functions



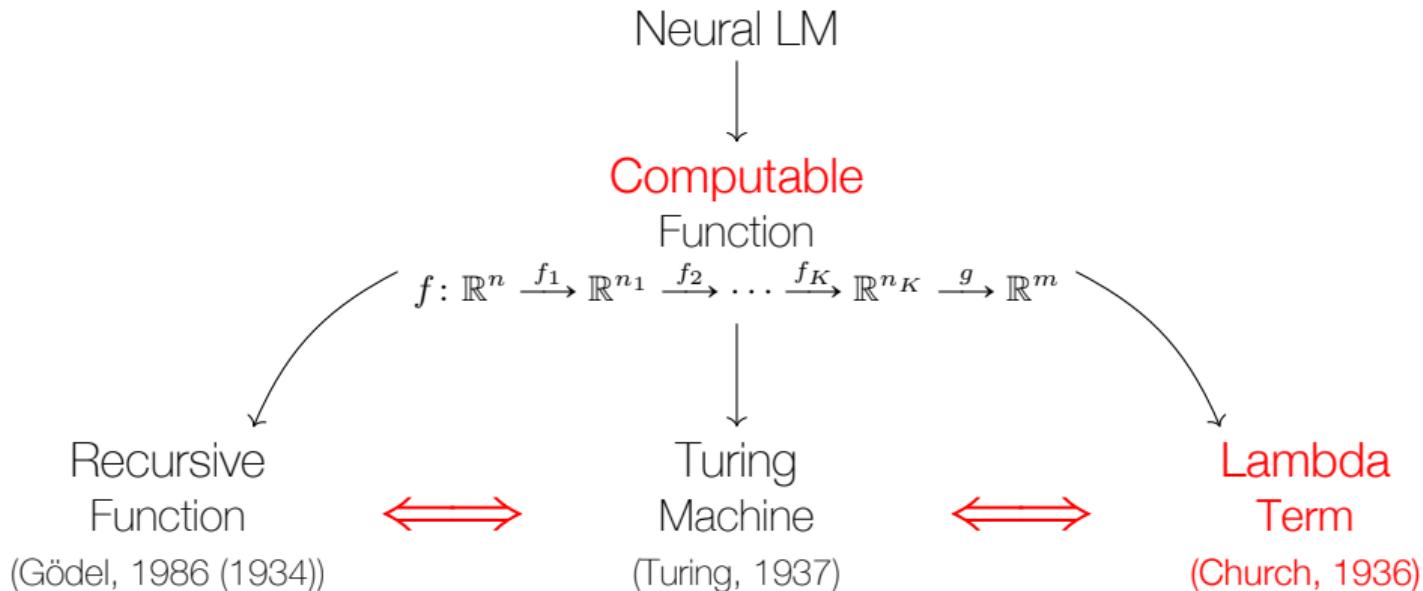
Neural LMs as Computable Functions



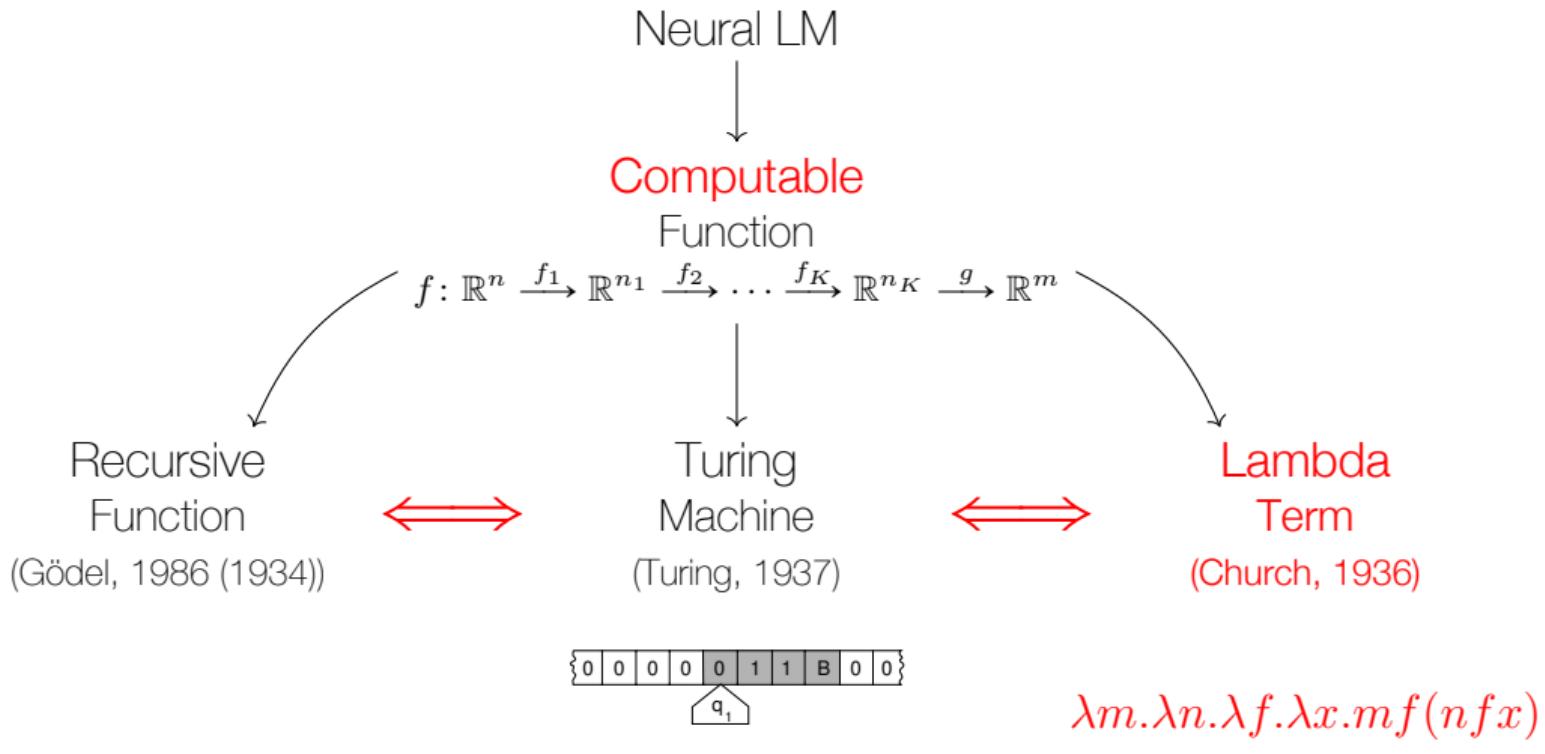
Neural LMs as Computable Functions



Neural LMs as Computable Functions



Neural LMs as Computable Functions



credit: Nynexman4464

Empirical Evaluation

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

Empirical Evaluation

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

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

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

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

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

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

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

...

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

Empirical Evaluation

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

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

$\lambda m. \lambda n. \lambda f. \lambda x. mf(nfx) (\lambda f. \lambda x. f(fx)) (\lambda f. \lambda x. f(f(fx)))$

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

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

...

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

Empirical Evaluation

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

0:	$\lambda f. \lambda x. x$	$\lambda m. \lambda n. \lambda f. \lambda x. mf(nfx)$
1:	$\lambda f. \lambda x. fx$	$(\lambda f. \lambda x. f(fx))$
2:	$\lambda f. \lambda x. f(fx)$	$(\lambda f. \lambda x. f(f(fx)))$
3:	$\lambda f. \lambda x. f(f(fx)))$	\vdots
4:	$\lambda f. \lambda x. f(f(f(fx))))$	\vdots
5:	$\lambda f. \lambda x. f(f(f(f(fx))))$	\vdots
...		\vdots
$n:$	$\lambda f. \lambda x. \underbrace{f(\dots(f}_{n \text{ times}} x) \dots)$	$\lambda f. \lambda x. f(f(f(f(f(fx))))))$

Empirical Evaluation

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

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

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

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

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

↓

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

↓

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

↓

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

↓

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

↓

...

↓

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

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

Interpretability

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

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

Interpretability

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

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

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

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

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

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

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

...

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

$$\lambda m. \lambda n. \lambda f. \lambda x. mf(nfx)(\lambda f. \lambda x. f(fx))(\lambda f. \lambda x. f(f(fx)))$$

$$\lambda m. \lambda n. \lambda f. \lambda x. mf(nfx)(\lambda g. \lambda y. g(gy))(\lambda h. \lambda z. h(h(hz)))$$

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

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

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

$$\lambda f. \lambda x. (\lambda y. f(fy))((\lambda h. \lambda z. h(h(hz)))fx)$$

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

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

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

Interpretability

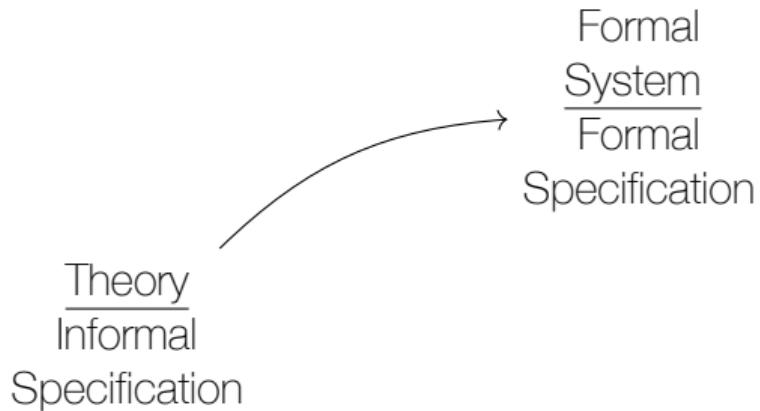
$P := \lambda m. \lambda n. \lambda f. \lambda x. mf(nfx)$

$P'' := \lambda RofAOe\tilde{N}5E | Ax\tilde{n}=\infty u \text{ ýmWf286ey'SOú>v&ì} \rightarrow 2 \text{ óÉ7öçoo}\{\tilde{a}>2flB^{\circ}\mu G\#\tilde{A}9\zeta U$
 $\infty btYBôY \text{ Ù } \ddot{e}\%_3;5 \text{ å[l-} \dot{e}u\ddot{o} \text{ Ü } 7-\ddot{U}. \lambda: \tilde{4m\tilde{O}\tilde{O}Y} \text{ 'è} \rightarrow \tilde{Is\ddot{O},\$+g\ddot{i},B^{\text{TM}}\div o-\#i\ddot{Y}\ddot{e} \text{ Üv}$
 $-g\ddot{O}\ddot{y}/\ddot{e}ijO\ddot{t}\ddot{C}fi \bullet J1«\ddot{E}\ddot{\phi},\ddot{I} \text{ h\ddot{a}et\ddot{t}\ddot{æ}Y\$^6 FiW»R\ddot{U}Kg\ddot{e} \text{ '}. \lambda\ddot{t}d^- \dots D2\div \ddot{o} \text{ ' x\ddot{e}\ddot{E}y. } \ddot{O} \text{ 'cb}$
 $Bé\ddot{E}N\ddot{E}1\ddot{E}\ddot{f}/\ddot{U}9\ddot{N}\mu-/JY\ddot{C}\ddot{o}\ddot{E}9\ddot{y}\ddot{A}\ddot{E}. \lambda\ddot{A}\ddot{I} \text{ '}\ddot{o}\ddot{C}, »fq\infty\pm\tilde{1}^B5\tilde{I}>O\tilde{g}^{\text{TM}}\tilde{6}\Omega e\tilde{a}\ddot{e}C/\tilde{a} \dots \ddot{O}$
 $\cdot f\ddot{O} \text{ '}\ddot{A}]\ddot{N}\ddot{a}y\ddot{E}\ddot{N}^\circ\ddot{E} \text{ '}. \lambda\ddot{E}\ddot{a}\ddot{e}\ddot{f}U\ddot{o}fE\ddot{U}\ddot{I} \text{ 'm\#,,4\sqrt{-}\div}\ddot{I}\ddot{p}\ddot{o} »y\ast v\ddot{t}\ddot{A}\ddot{J}\ddot{A}\ddot{F}\ddot{1}\ddot{u}\ddot{A}\ddot{o}\ddot{z}\ddot{<}\ddot{n}\ddot{M}\ddot{D}\ddot{j}\ddot{C}\ddot{E}$
 $B\ddot{E}\ddot{e}\ddot{I}\ddot{T} \text{ '}\ddot{E}\ddot{a}\%_0\ddot{A}\ddot{C}\ddot{\Omega} @\ddot{[\ddot{\varnothing}\ddot{^}\ddot{~}]} \ddot{I}\ddot{h}\ddot{t}: \tilde{4m\tilde{O}\tilde{O}Y} \text{ 'è} \rightarrow \tilde{Is\ddot{O},\$+g\ddot{i},B^{\text{TM}}\div o-\#i\ddot{Y}\ddot{e} \text{ Üv-g\ddot{O}\ddot{y}}$
 $/\ddot{e}ijO\ddot{t}\ddot{C}fi \bullet J1«\ddot{E}\ddot{\phi},\ddot{I} \text{ h\ddot{a}et\ddot{t}\ddot{æ}Y\$^6 FiW»R\ddot{U}Kg\ddot{e} \text{ '}\ddot{A}\ddot{I} \text{ '}\ddot{o}\ddot{C}, »fq\infty\pm\tilde{1}^B5\tilde{I}>O\tilde{g}^{\text{TM}}\tilde{6}$
 $\Omega e\tilde{a}\ddot{e}C/\tilde{a} \dots \ddot{O} \cdot f\ddot{O} \text{ '}\ddot{A}]\ddot{N}\ddot{a}y\ddot{E}\ddot{N}^\circ\ddot{E} \text{ '}(\ddot{t}d^- \dots D2\div \ddot{o} \text{ ' x\ddot{e}\ddot{E}y. } \ddot{O} \text{ 'cbBé\ddot{E}N\ddot{E}1\ddot{E}\ddot{f}/\ddot{U}9\ddot{N}\mu-/}$
 $JY\ddot{C}\ddot{o}\ddot{E}9\ddot{y}\ddot{A}\ddot{E}\ddot{A}\ddot{I} \text{ '}\ddot{o}\ddot{C}, »fq\infty\pm\tilde{1}^B5\tilde{I}>O\tilde{g}^{\text{TM}}\tilde{6}\Omega e\tilde{a}\ddot{e}C/\tilde{a} \dots \ddot{O} \cdot f\ddot{O} \text{ '}\ddot{A}]\ddot{N}\ddot{a}y\ddot{E}\ddot{N}^\circ\ddot{E} \text{ '}\ddot{E}\ddot{a}\ddot{e}\ddot{f}U\ddot{o}fE\ddot{U}\ddot{I} \text{ 'm\#,,4\sqrt{-}\div}\ddot{I}\ddot{p}\ddot{o} »y\ast v\ddot{t}\ddot{A}\ddot{J}\ddot{A}\ddot{F}\ddot{1}\ddot{u}\ddot{A}\ddot{o}\ddot{z}\ddot{<}\ddot{n}\ddot{M}\ddot{D}\ddot{j}\ddot{C}\ddot{E}\ddot{B}\ddot{E}\ddot{e}\ddot{I}\ddot{T} \text{ '}\ddot{E}\ddot{a}\%_0\ddot{A}\ddot{C}\ddot{\Omega} @\ddot{[\ddot{\varnothing}\ddot{^}\ddot{~}]} \ddot{I}\ddot{h}\ddot{t})(\ddot{E}\ddot{I}\ddot{U}\ddot{e}\ddot{í}\ddot{4}\ddot{W}\ddot{\mu}\ddot{I} \text{ '}\ddot{w},\ddot{\$}\ddot{\Omega}\ddot{“}\ddot{K}\ddot{5}\ddot{e}\ddot{A}\ddot{\P}\ddot{3}[m \text{ '}\ddot{B}\ddot{A}\ddot{f}\ddot{f}\ddot{O}; \ddot{o}\ddot{J}\ddot{c}\ddot{C}\ddot{E}\ddot{\tilde{o}}\ddot{Y}\ddot{O}\ddot{c}\ddot{B},\ddot{\$}\ddot{\tilde{A}}\ddot{a}\ddot{a}\ddot{}}\ddot{O}\ddot{A}\ddot{\tilde{O}}\ddot{3};$
 $\ddot{?}\ddot{o}\ddot{-}\ddot{o}\ddot{C}\ddot{E}\ddot{@}\ddot{f}\ddot{8}\ddot{R}\ddot{C}\ddot{æ}\ddot{e}\ddot{o}\ddot{*}\ddot{\<}\ddot{Y}\ddot{-}\ddot{o}\ddot{1}\ddot{2}\ddot{A}\ddot{\%}\ddot{0}\ddot{a}\ddot{O}\ddot{Ü}\ddot{\#}\ddot{i}\ddot{",}\ddot{u}\ddot{"}\ddot{\<}\ddot{\hat{o}},\ddot{\infty}\ddot{\hat{I}\ddot{a}\ddot{a}\ddot{a}\ddot{}}\ddot{\o}\ddot{A}\ddot{d}\ddot{|}\ddot{\tilde{N}}\ddot{'}\ddot{E}\ddot{y}\ddot{\varnothing};\ddot{^}\ddot{W}$
 $\ddot{\>}\ddot{w}\ddot{o}[\ddot{\}]\ddot{\>}\ddot{Ö}\ddot{E}\ddot{u}\ddot{w}\ddot{'}\ddot{6}\ddot{<}\ddot{u}\ddot{^}\ddot{=}\ddot{a}\ddot{O}\ddot{-}\ddot{I}\ddot{D}\ddot{z}\ddot{?}\ddot{2}\ddot{\pm}\ddot{|}\ddot{é}\ddot{'}\ddot{3}\ddot{A}\ddot{/}\ddot{r}\ddot{x}\ddot{\mu}\ddot{\infty}\ddot{\mu}\ddot{\$}\ddot{\tilde{A}\ddot{e}\ddot{A}\ddot{*}\ddot{l}\ddot{f}\ddot{\sim}\ddot{\tilde{u}}}\ddot{'+}\ddot{I}\ddot{V}\ddot{iy}\ddot{^}\ddot{a}\ddot{G}\ddot{æ}\ddot{ß}\ddot{ä}\ddot{g}\ddot{o}\ddot{/},\ddot{u}\ddot{N}\ddot{)}$

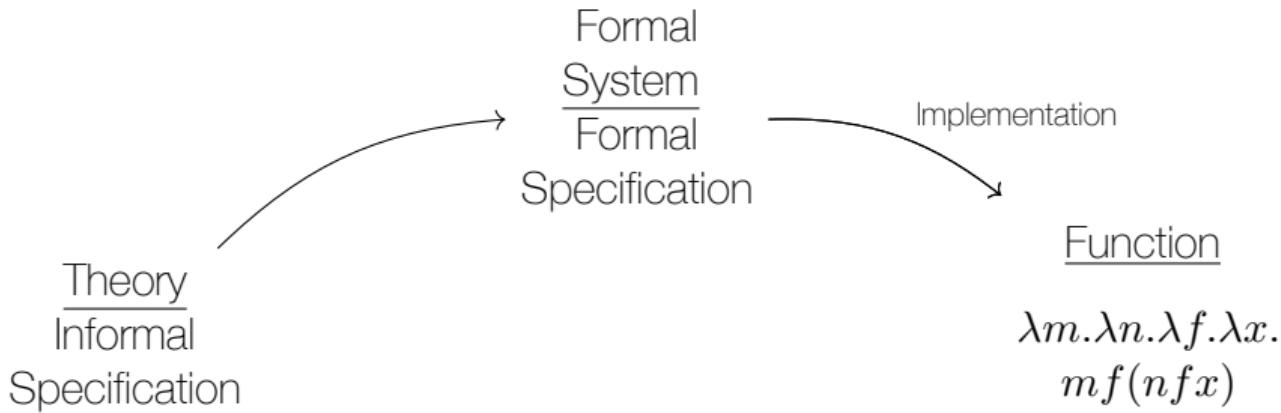
Making it Explicit

Theory
Informal
Specification

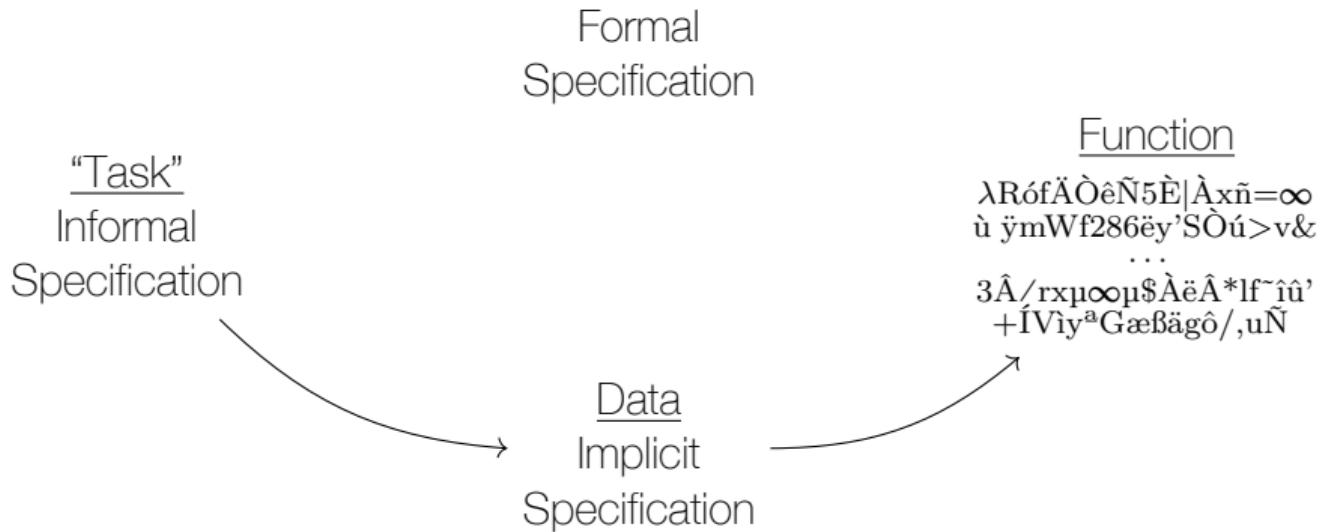
Making it Explicit



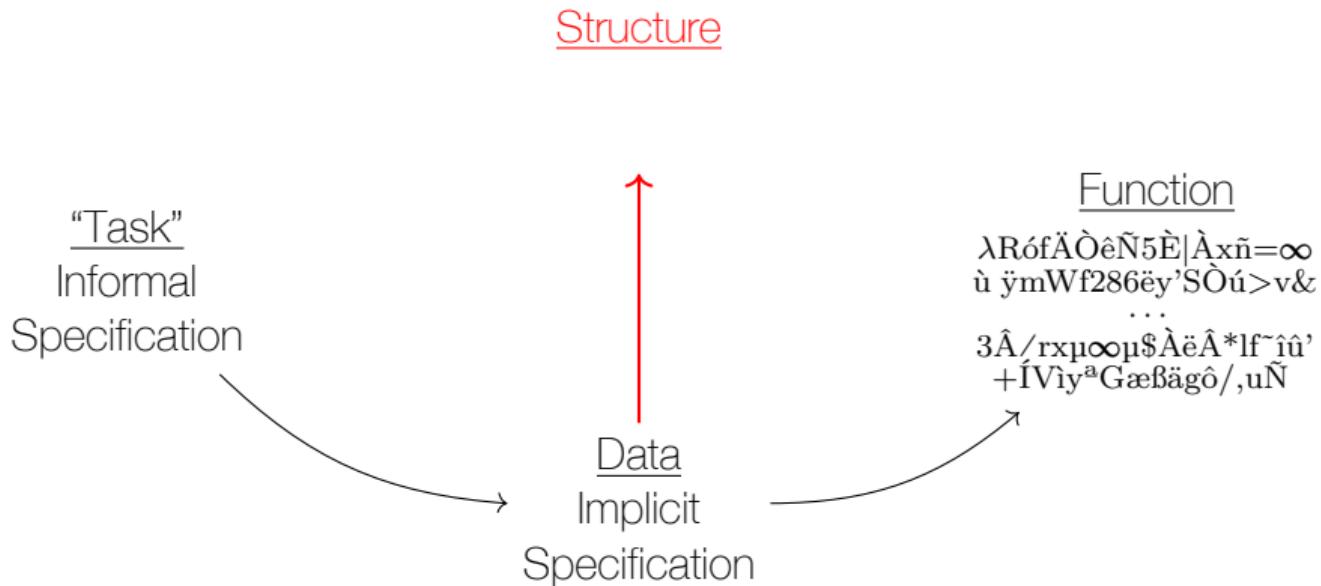
Making it Explicit



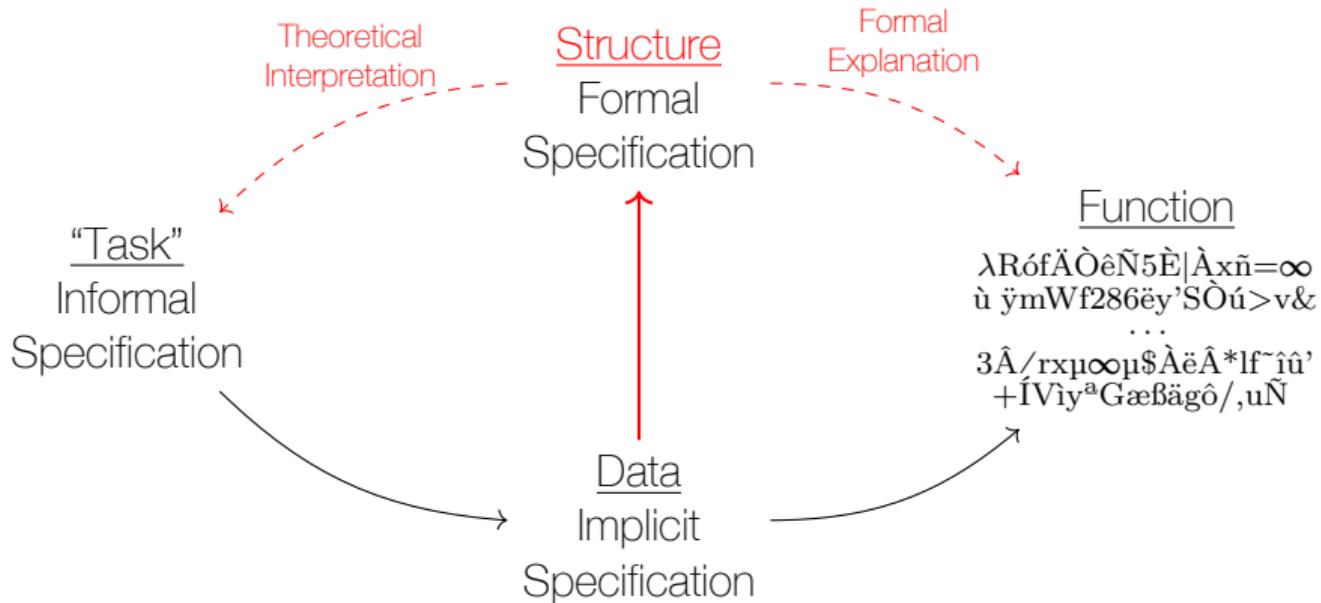
Making it Explicit



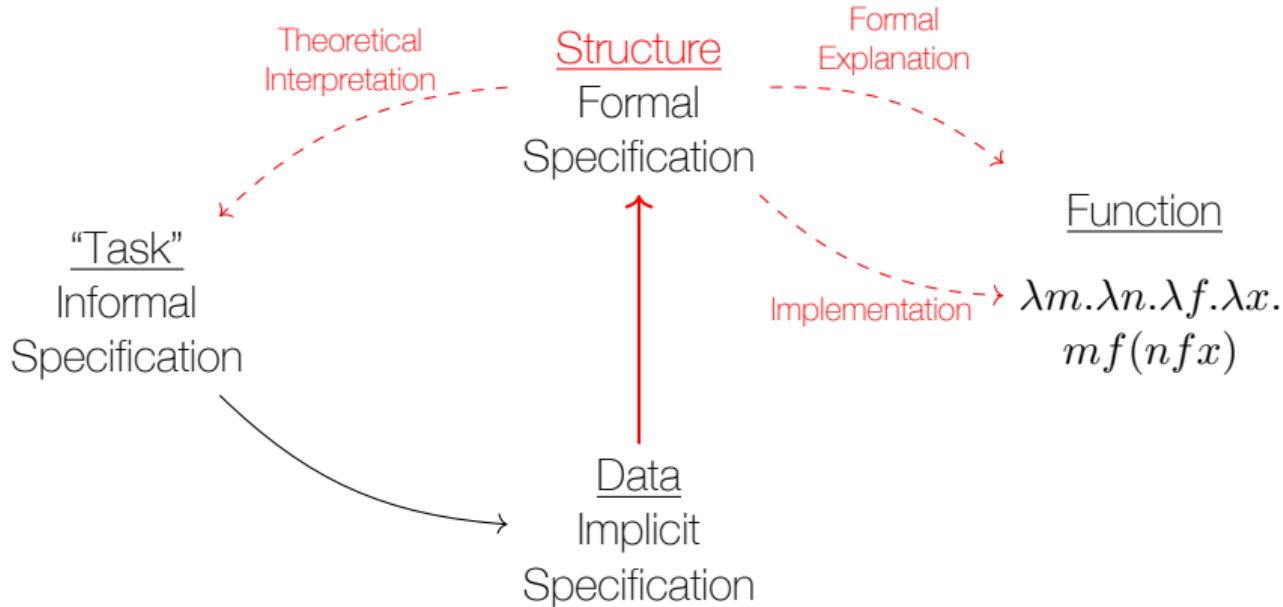
Making it Explicit



Making it Explicit



Making it Explicit



Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

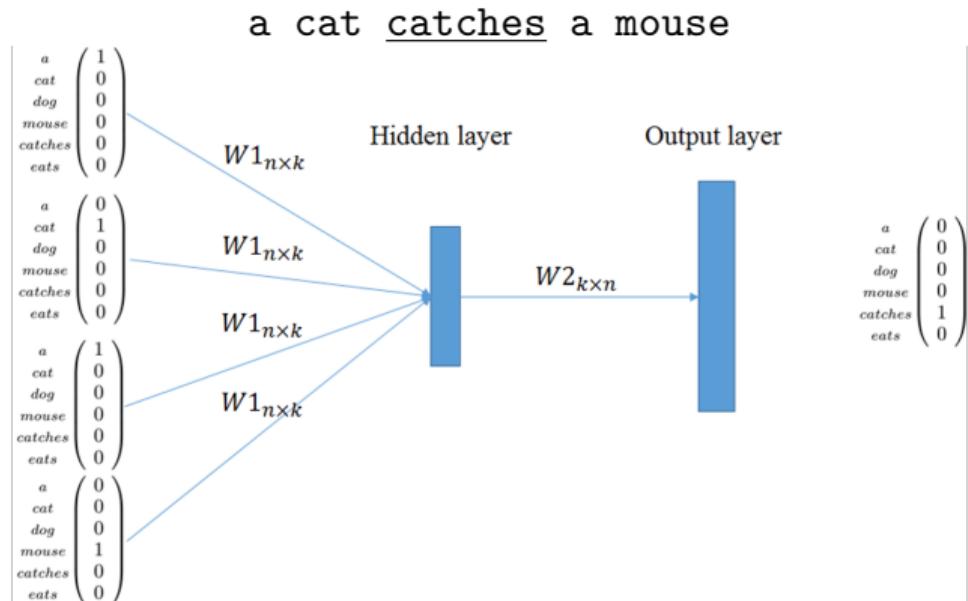
The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Machine-Learning the Embedding Space



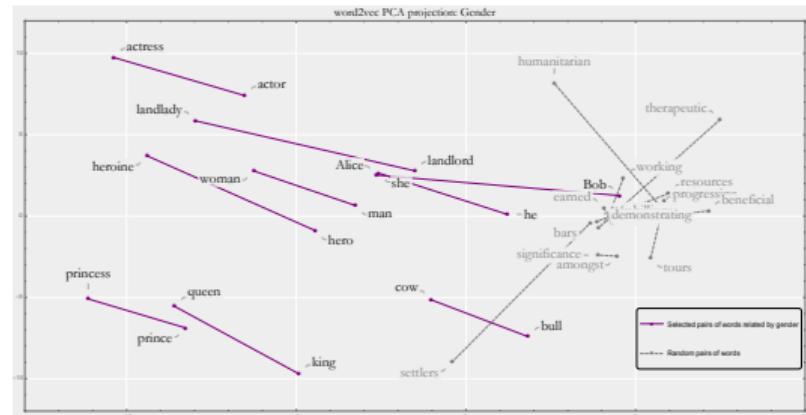
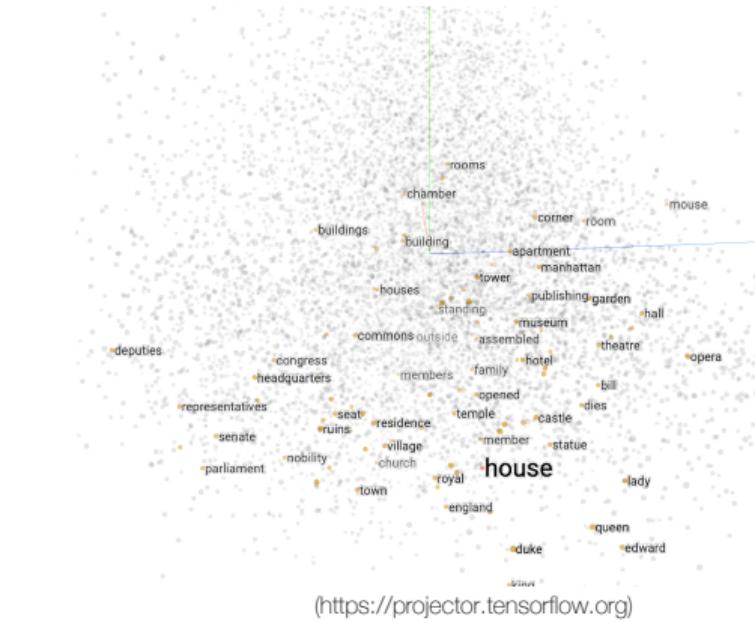
Credit: Ferrone et al., 2017

Credit: xkcd.com

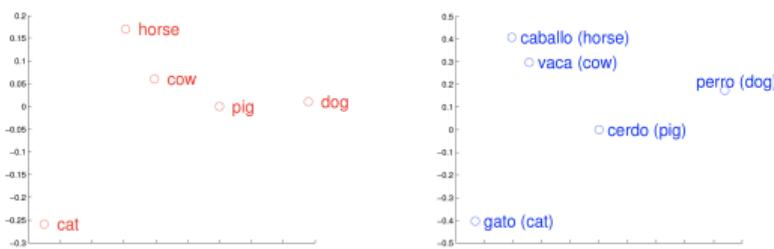
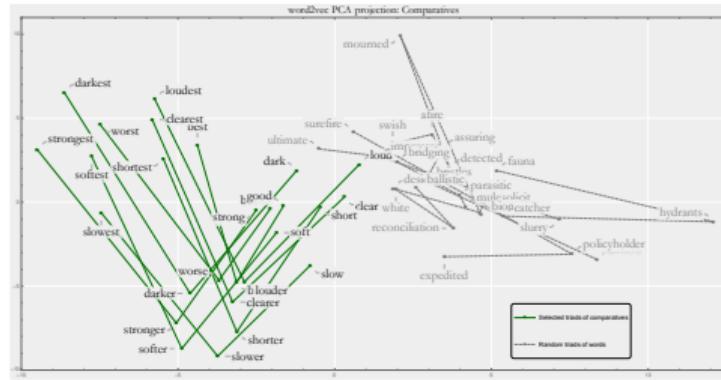
Juan Luis Gastaldi | The Structure, Not the Prompt

9/35

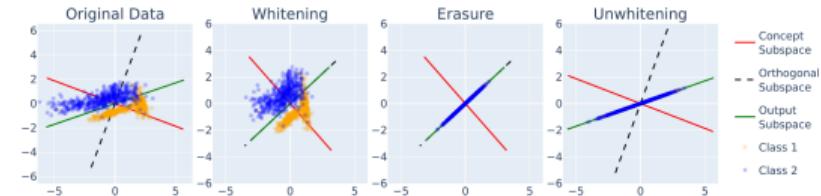
Embedding Space: Similarity and Analogy



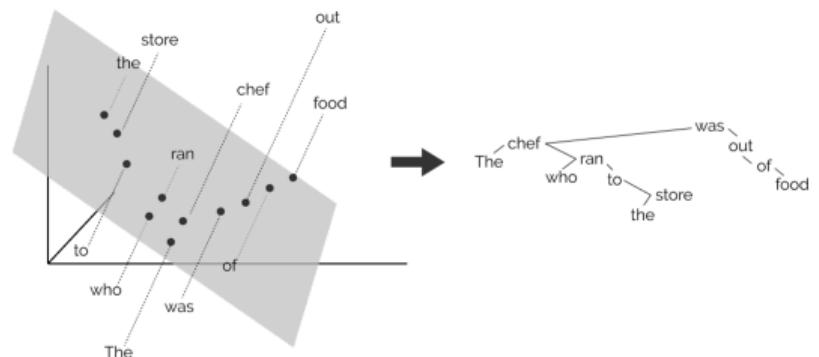
Embedding Space: Other Applications



(Mikolov et al., 2013)



(Belrose et al., 2024)



(<https://nlp.stanford.edu/~johnhew/structural-probe.html>)

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)

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

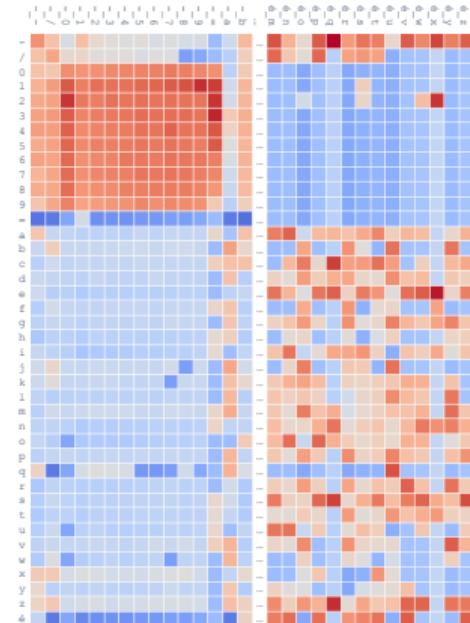
$$\frac{\partial \ell}{\partial (\vec{w} \cdot \vec{c})} = 0 \quad \text{when} \quad \vec{w} \cdot \vec{c} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

- 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

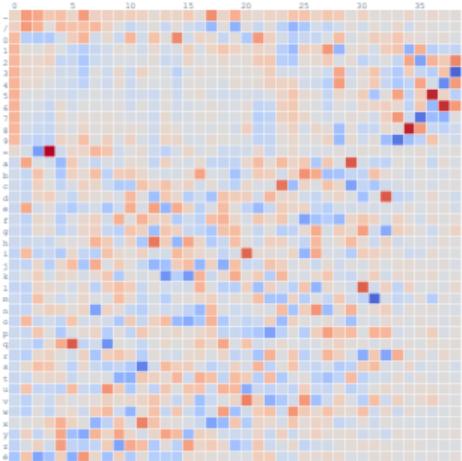
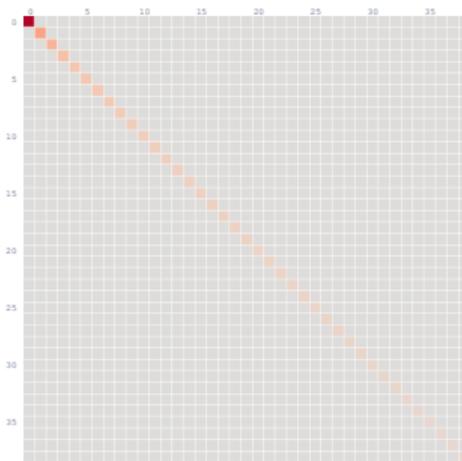
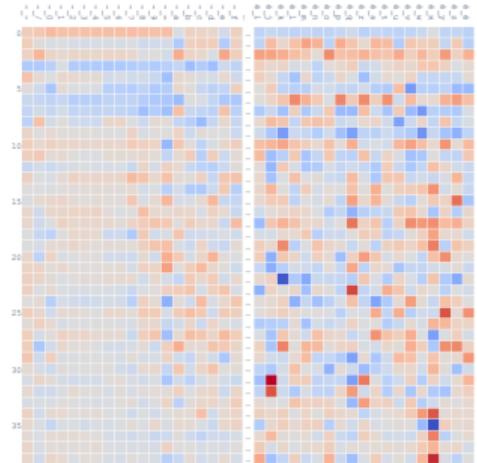
$$W = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, \dots, w, x, y, z, é\}$$

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



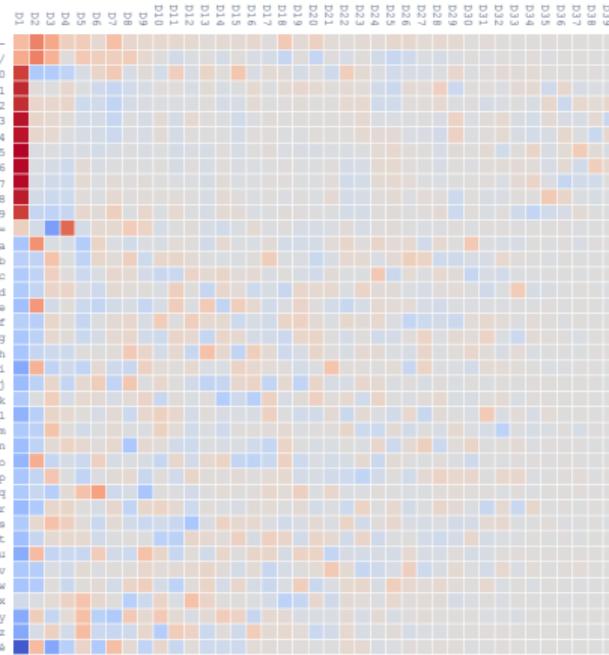
$$\begin{aligned} M_{wc} &= \text{pmi}(w, c) \\ &= \log \frac{p(w, c)}{p(w)p(c)} \end{aligned}$$

SVD of Wikipedia Character PMI Matrix

 U  Σ  V^T 

Truncate

$U \times \Sigma$

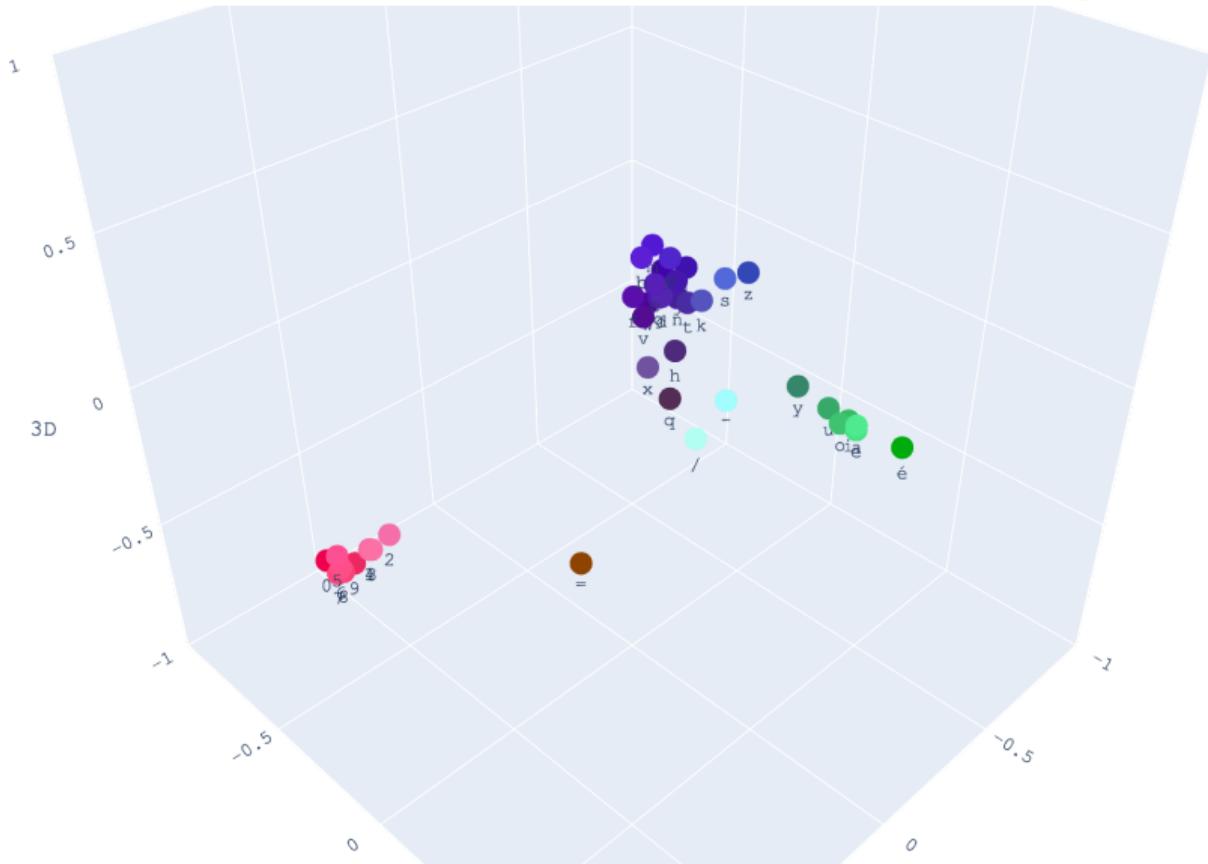


Truncate

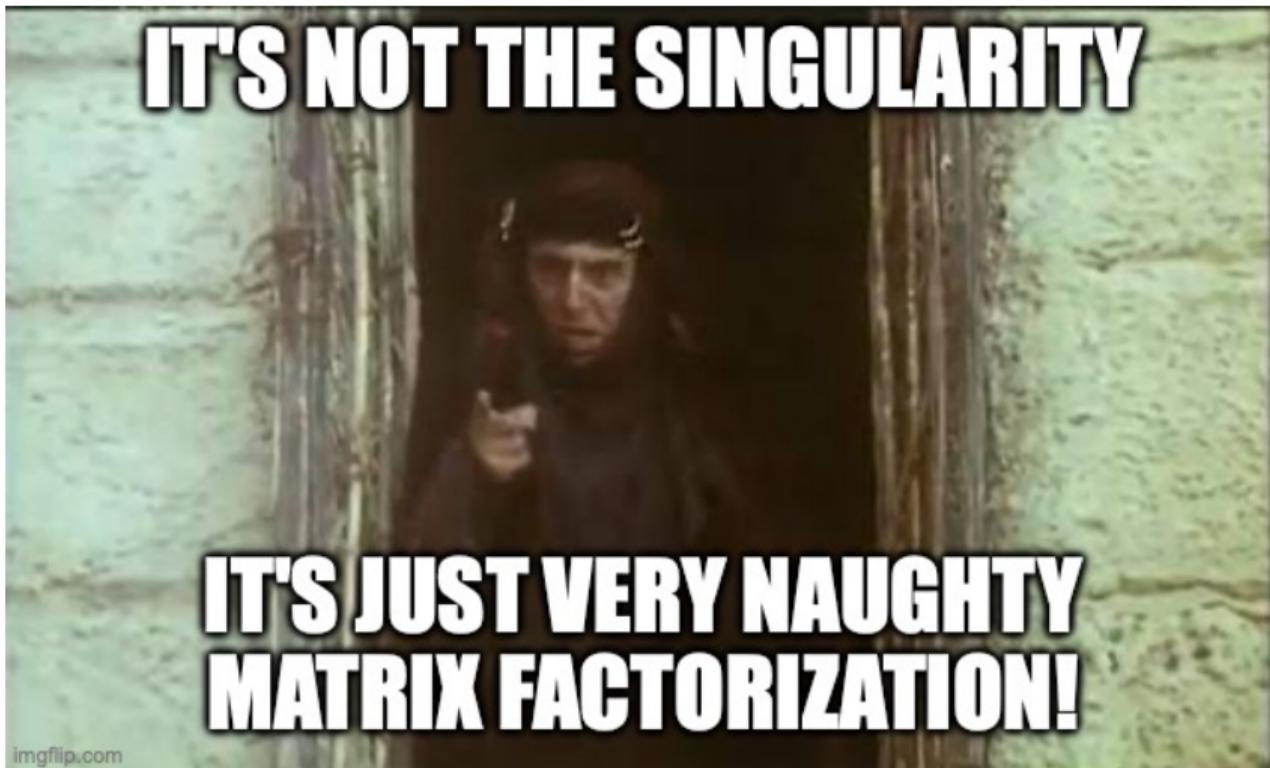
$\hat{U} \times \hat{\Sigma}$



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



What to conclude?



But Why?

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)

Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Embeddings as Functions Over Sets

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

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

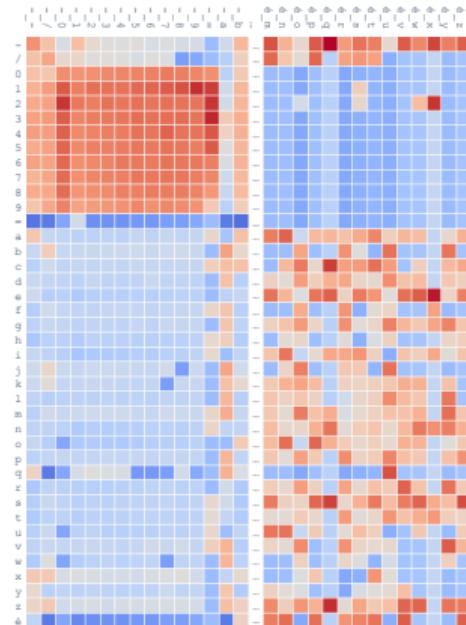
Embeddings as Functions Over Sets

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

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

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

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



Embeddings as Functions Over Sets

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

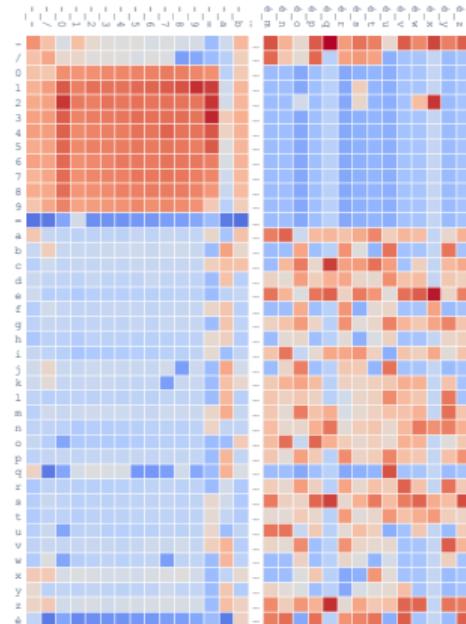
$$\textcolor{blue}{Y} = X \times X = \{(-, -), (-, /), (-, 0), \dots, (\text{é}, \text{z}), (\text{é}, \text{é})\}$$

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

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

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

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



Embeddings as Functions Over Sets

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

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

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

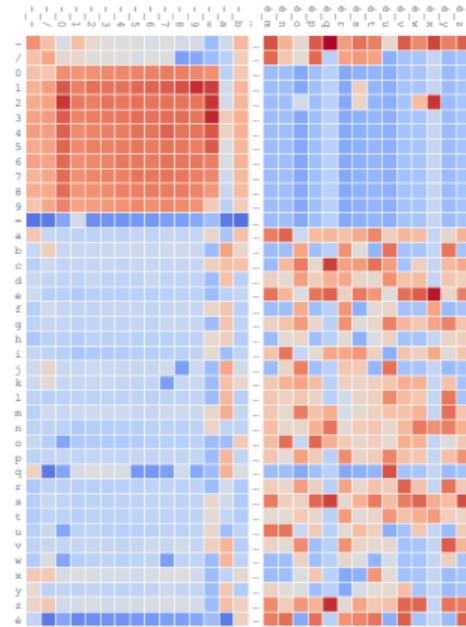
$$(\textcolor{red}{x}, \textcolor{blue}{y}) \mapsto \text{pmi}(\textcolor{red}{x}, \textcolor{blue}{y})$$

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

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

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

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



Embeddings as Functions Over Sets

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

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

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

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

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

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

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

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

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

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

Embeddings as Functions Over Sets

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

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

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

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

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

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

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

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

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

Embeddings as Functions Over Sets

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

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

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

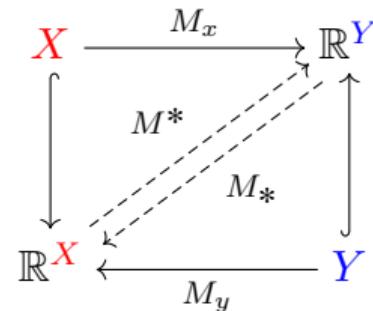
$$(\textcolor{red}{x}, \textcolor{blue}{y}) \mapsto \text{pmi}(\textcolor{red}{x}, \textcolor{blue}{y})$$

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

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

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

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

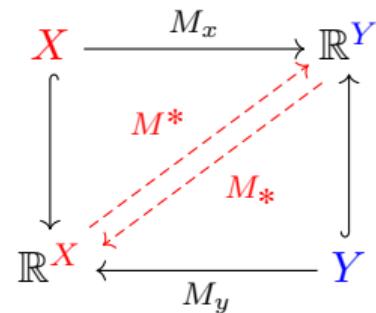


$$M^*: \mathbb{R}^{\textcolor{red}{X}} \rightarrow \mathbb{R}^{\textcolor{blue}{Y}}$$

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

Embeddings as Functions Over Sets

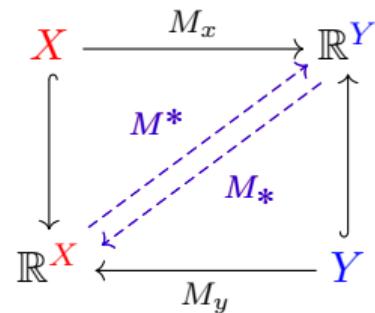
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$



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

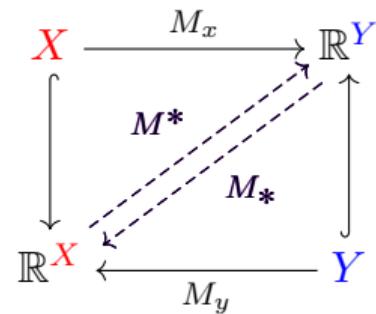
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

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

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

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

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



Embeddings as Functions Over Sets

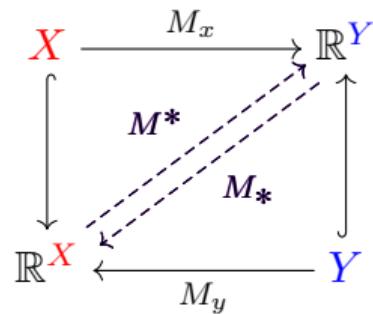
$$M_* M^* : \mathbb{R}^X \rightarrow \mathbb{R}^X$$

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

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

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

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



$$U := [\mathbf{u}_1, \dots, \mathbf{u}_m]$$

$$M = U \Sigma V^T \quad V := [\mathbf{v}_1, \dots, \mathbf{v}_n]$$

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

Embeddings as Functions Over Sets

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

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

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

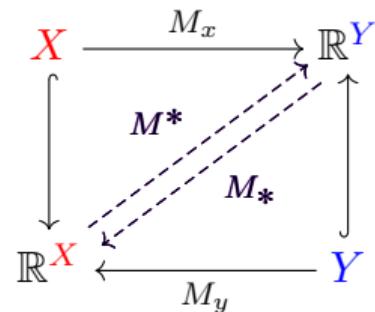
$$\{v_1, \dots, v_n\} \subset \mathbb{R}^Y$$

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

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

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

The u_i and v_i are (linear)
fixed points!

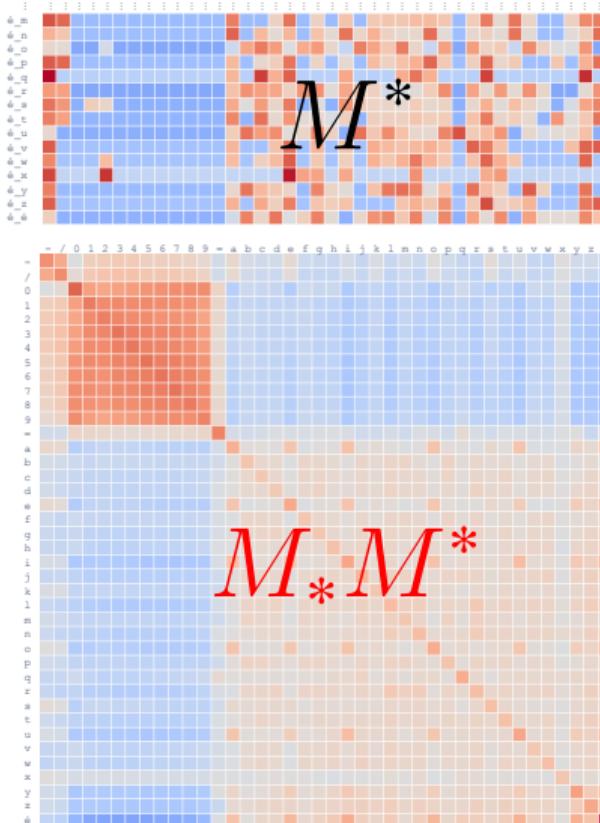
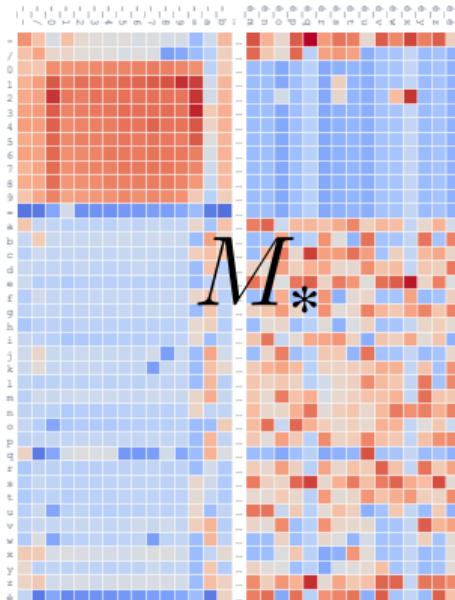


$$U := [\underline{u_1}, \dots, \underline{u_m}]$$

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

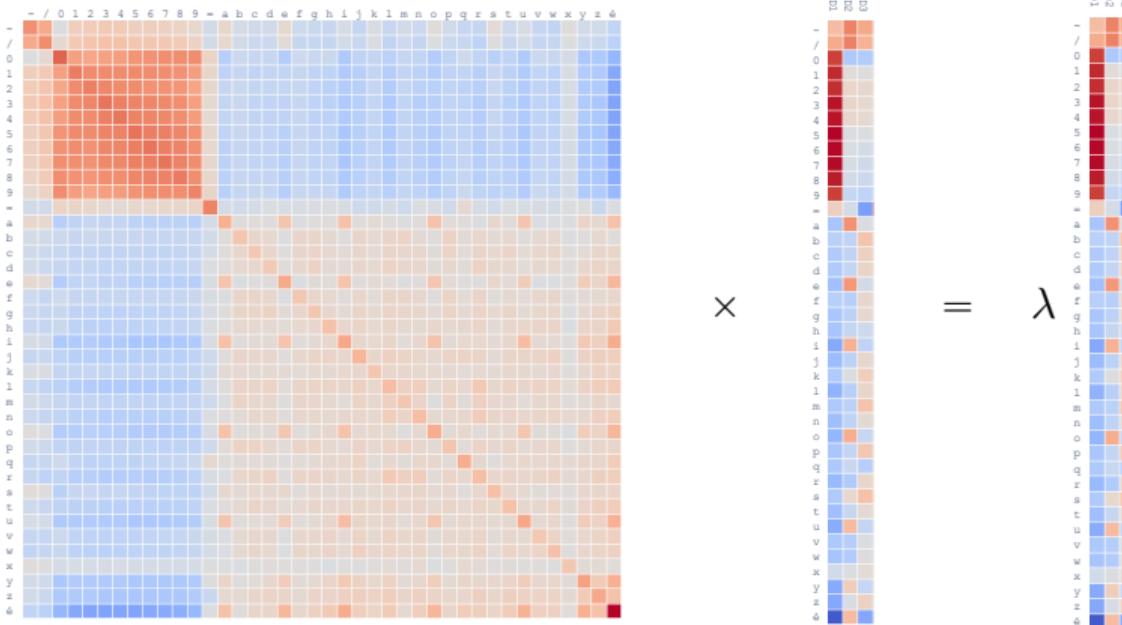
$$\Sigma := \begin{bmatrix} \sqrt{\lambda_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sqrt{\lambda_r} \end{bmatrix}$$

$M_* M^*$ as a Covariance Matrix



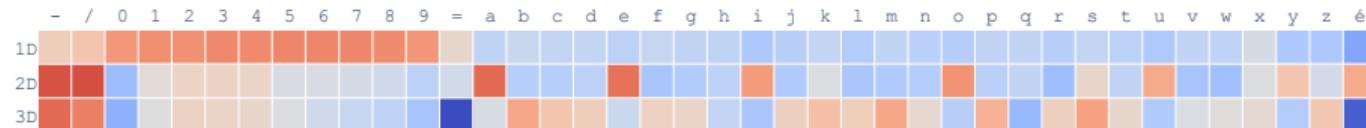
Eigenvectors as Fixed Points

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

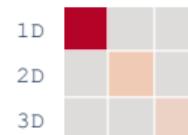


Structural Features

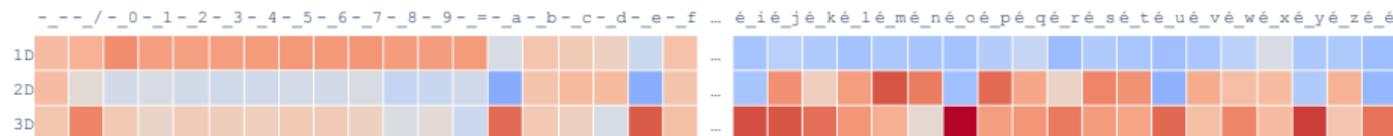
Eigenvectors of $M_* M^*$:



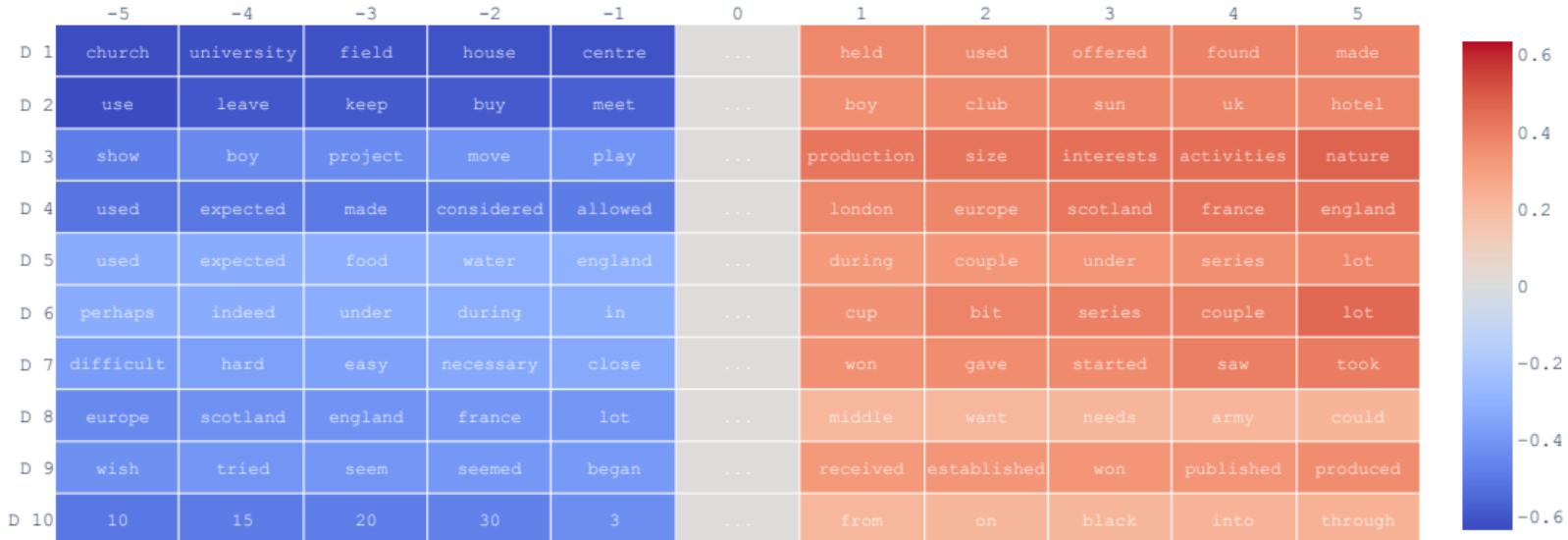
Eigenvalues of $M_* M^*$ and $M^* M_*$:



Eigenvectors of $M^* M_*$:



Words



Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

The Algebra Behind the Embeddings

The Structure Behind the Algebra

The Categories Behind the Structure

Conclusion

Embeddings as Functors Over Categories

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

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

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

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

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

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

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

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

$$\begin{array}{ccc} \textcolor{red}{X} & \xrightarrow{M_x} & \mathbb{R}^{\textcolor{blue}{Y}} \\ \downarrow & M^* \nearrow & \uparrow \\ \mathbb{R}^{\textcolor{red}{X}} & \xleftarrow[M_y]{\quad} & \textcolor{blue}{Y} \end{array}$$

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

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

Embeddings as Functors Over Categories

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

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

Profunctor

$$\mathcal{M}: \textcolor{orange}{C}^{\text{op}} \times \textcolor{blue}{D} \rightarrow \text{Set}$$

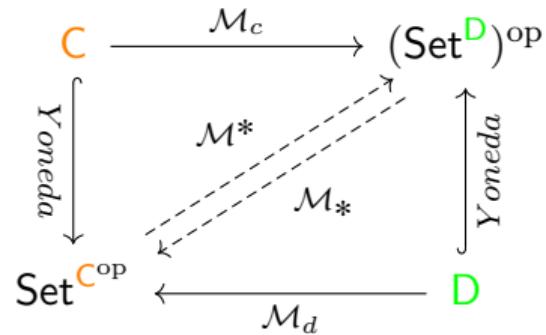
$$(\textcolor{orange}{c}, \textcolor{blue}{d}) \mapsto \mathcal{M}(\textcolor{orange}{c}, \textcolor{blue}{d})$$

$$\mathcal{M}_c: \textcolor{orange}{C} \rightarrow (\text{Set}^{\textcolor{blue}{D}})^{\text{op}}$$

$$\textcolor{orange}{c} \mapsto \mathcal{M}(\textcolor{orange}{c}, -)$$

$$\mathcal{M}_d: \textcolor{blue}{D} \rightarrow \text{Set}^{\textcolor{orange}{C}^{\text{op}}}$$

$$\textcolor{blue}{d} \mapsto \mathcal{M}(-, \textcolor{blue}{d})$$



$$\mathcal{M}^*: \text{Set}^{\textcolor{orange}{C}^{\text{op}}} \rightarrow (\text{Set}^{\textcolor{blue}{D}})^{\text{op}}$$

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

Embeddings as Functors Over Categories

Isbell Adjunction

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

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

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

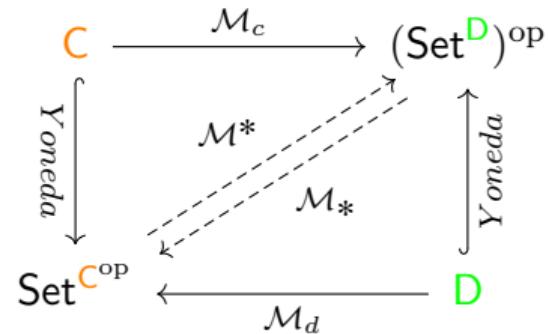
$$\text{Fix}(\mathcal{M}_* \mathcal{M}^*) := \{f \in \text{Set}^{\text{C}^{\text{op}}} \mid \mathcal{M}_* \mathcal{M}^*(f) \cong f\}$$

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

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

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

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



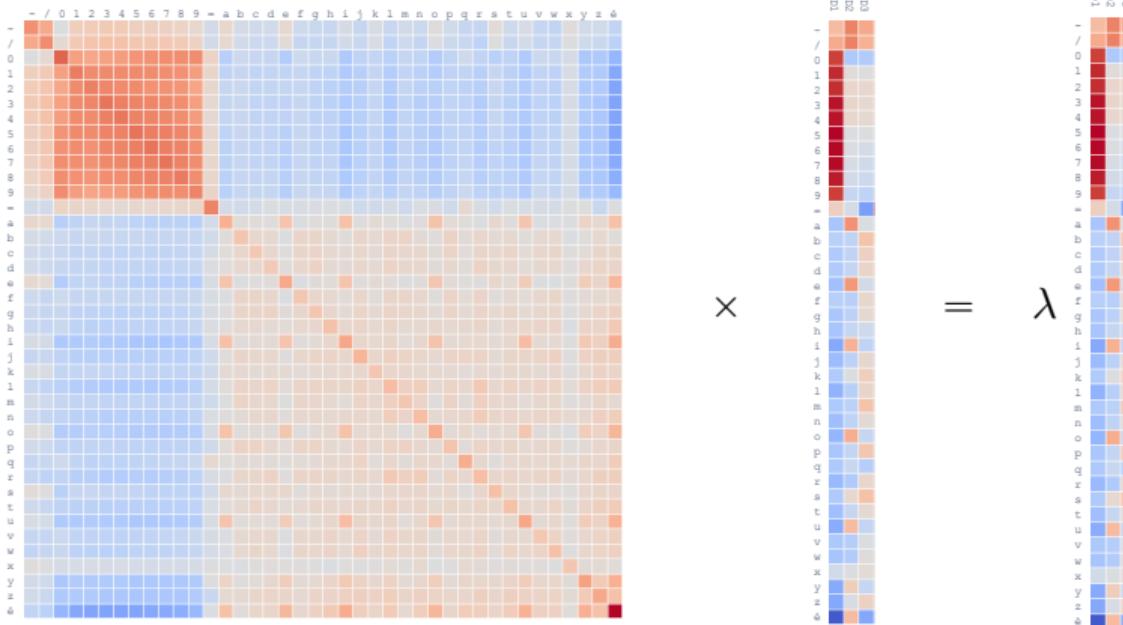
Categories **C** and **D**
can be enriched!

E.g.:

$$\begin{aligned} \mathcal{M}^*: 2^{\text{C}^{\text{op}}} &\leftrightarrows (2^{\text{D}})^{\text{op}}: \mathcal{M}_* \\ \mathcal{M}^*: \bar{\mathbb{R}}^{\text{C}^{\text{op}}} &\leftrightarrows (\bar{\mathbb{R}}^{\text{D}})^{\text{op}}: \mathcal{M}_* \end{aligned}$$

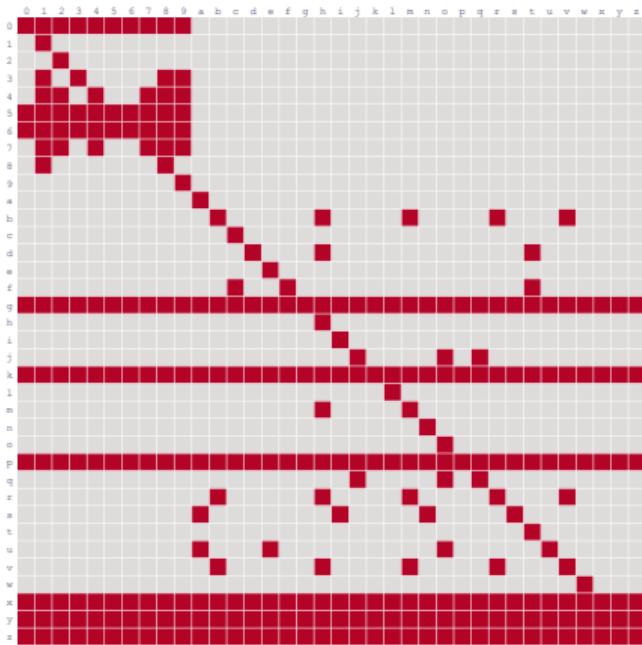
Binary Fixed Points

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



Binary Fixed Points

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



★



?

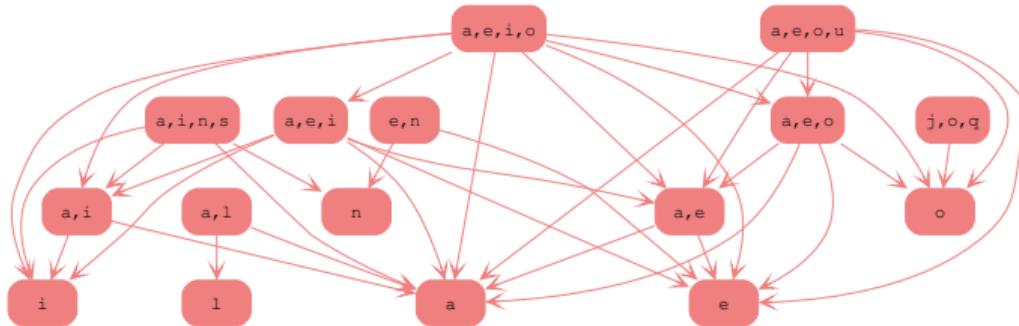
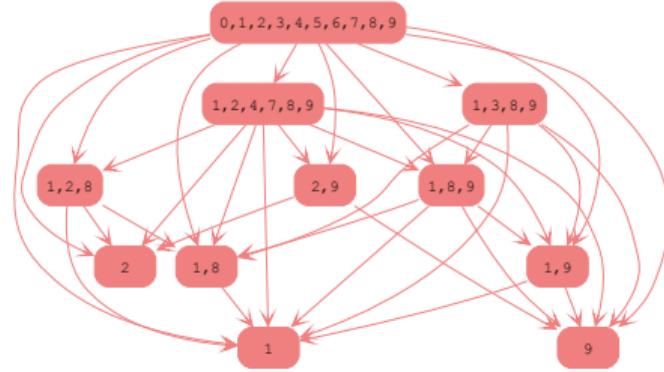
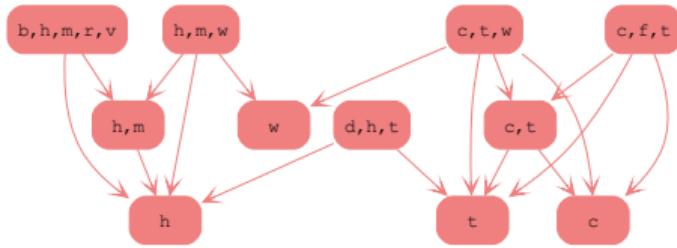


“Eigensets”

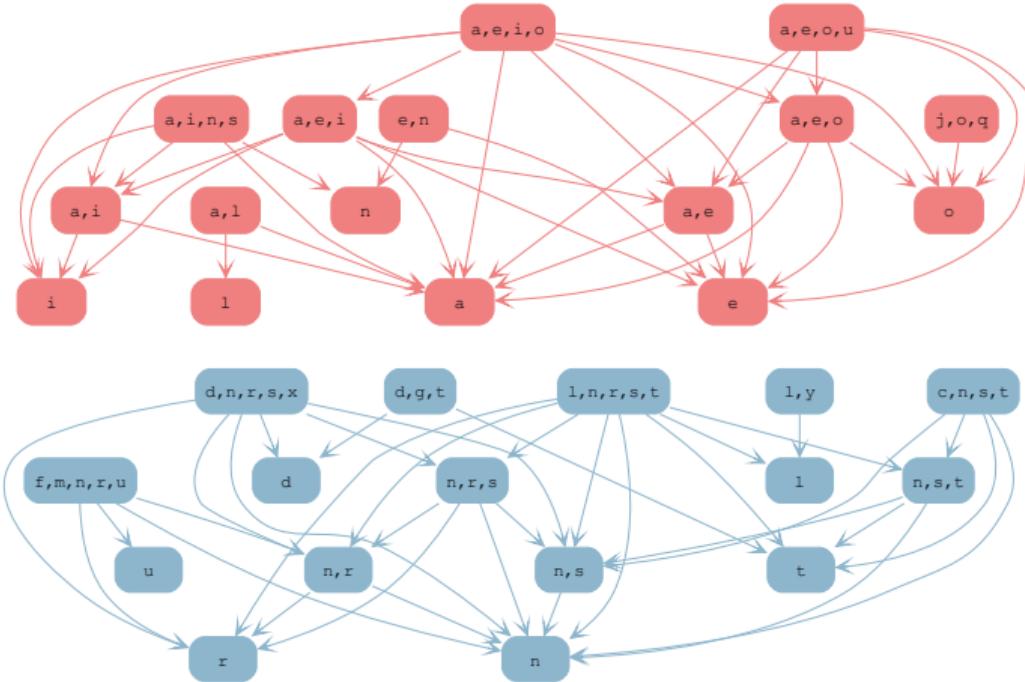
$$\mathcal{M}_*\mathcal{M}^*f = f$$

0,1,2,3,4,5,6,7,8,9	1,2,4,7,8,9	b,h,m,r,v	a,e,i,o	a,e,o,u	a,i,n,s	1,3,8,9
1,2,8	h,m,w	1,8,9	d,h,t	j,o,q	c,f,t	c,t,w
a,e,o	a,e,i	h,m	2,9	a,i	w	1,9
1,8	a,e	l	t	n	c	h
2	i	e	a	o	1	9
e,n	a,l	c,t				

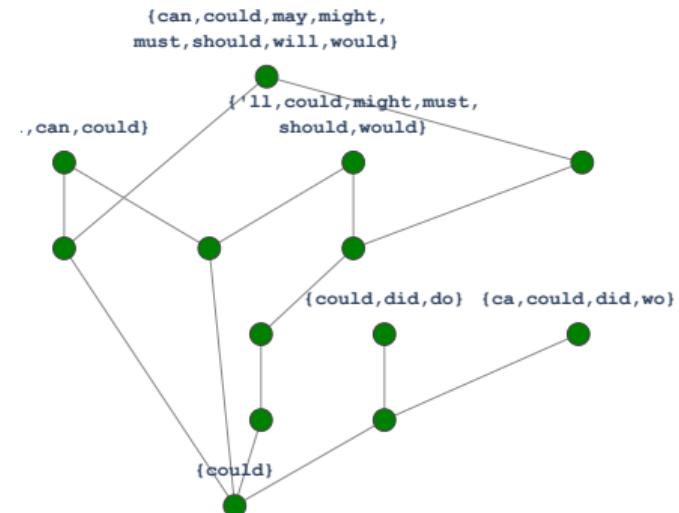
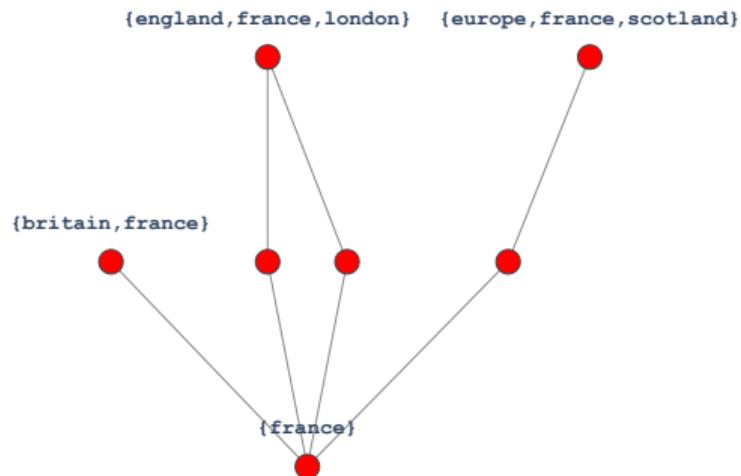
Which Structure?



Which Structure?

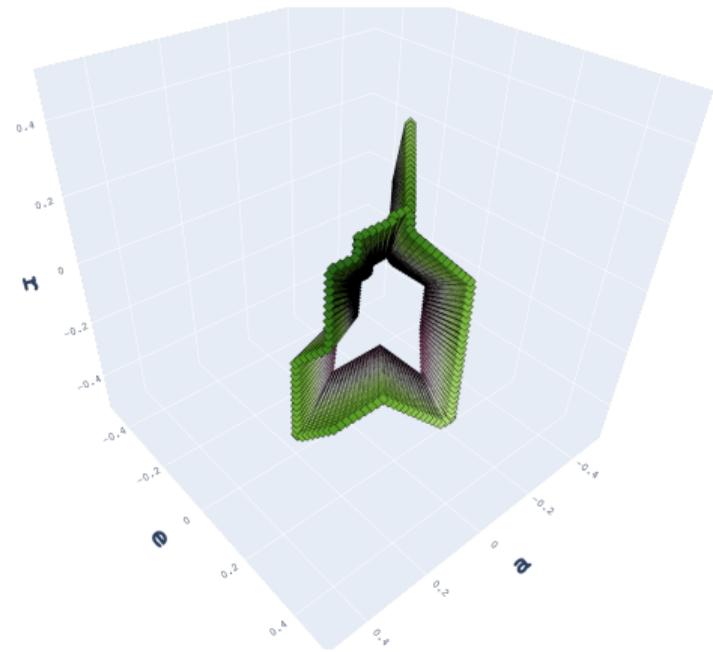
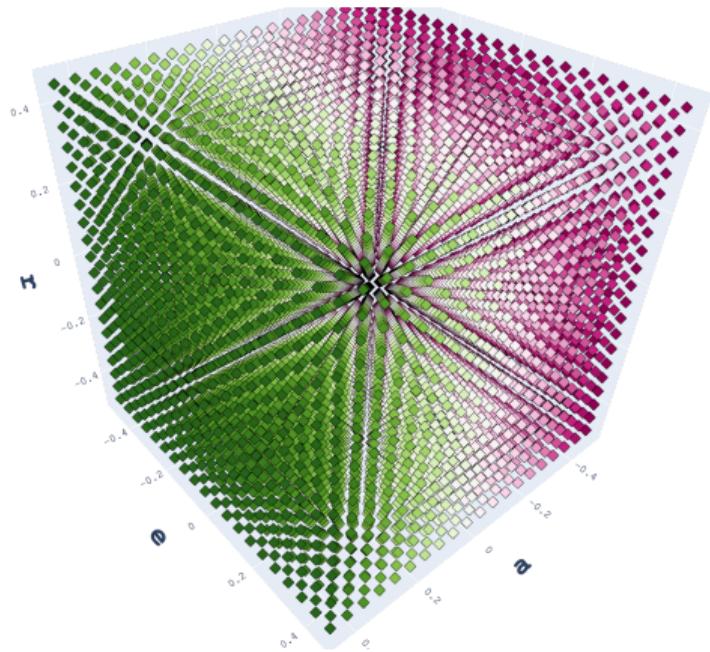


Formal Concepts (Words)

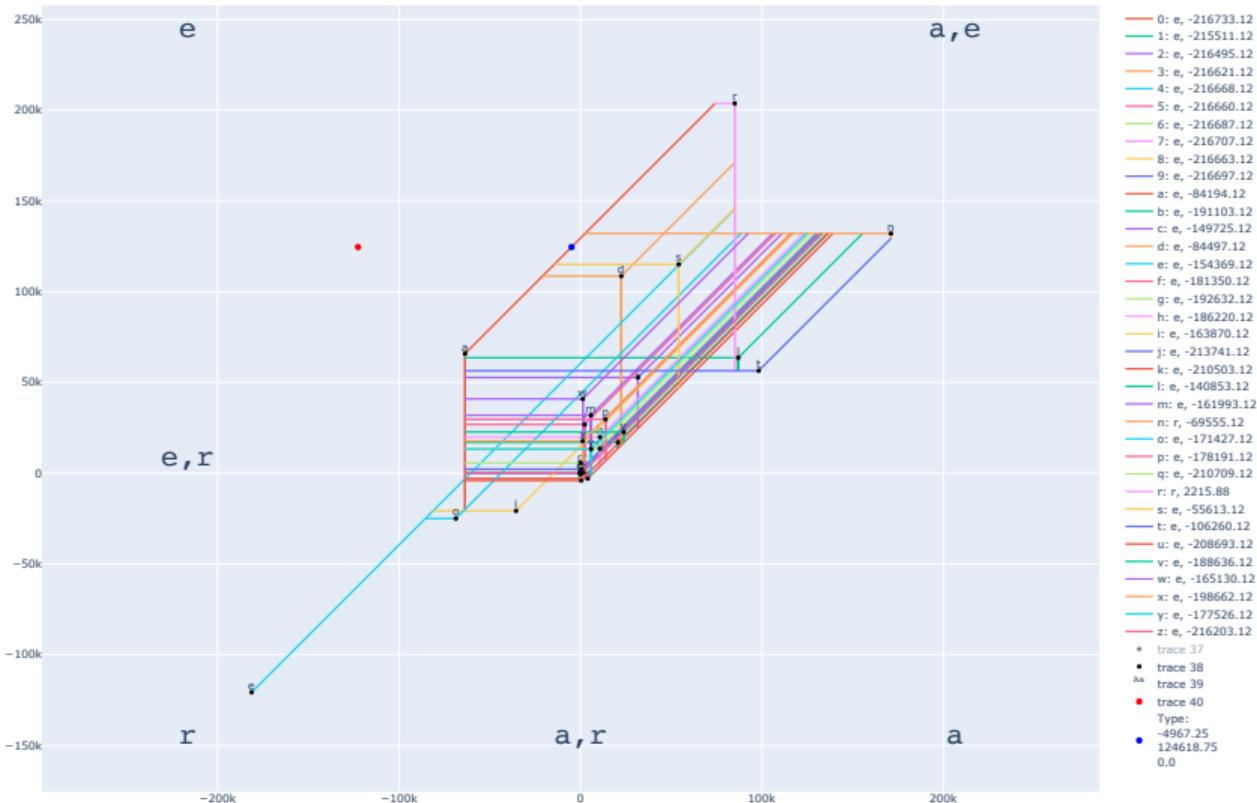


Nucleus

$$\bar{\mathbb{R}}^{\{a,e,r\}} \xrightarrow{\mathcal{M}_*\mathcal{M}^*} \bar{\mathbb{R}}^{\{a,e,r\}}$$



Internal Structure of the Nucleus



Theory of Computational Types

Definition (Polar/Orthogonal - Girard, 2011)

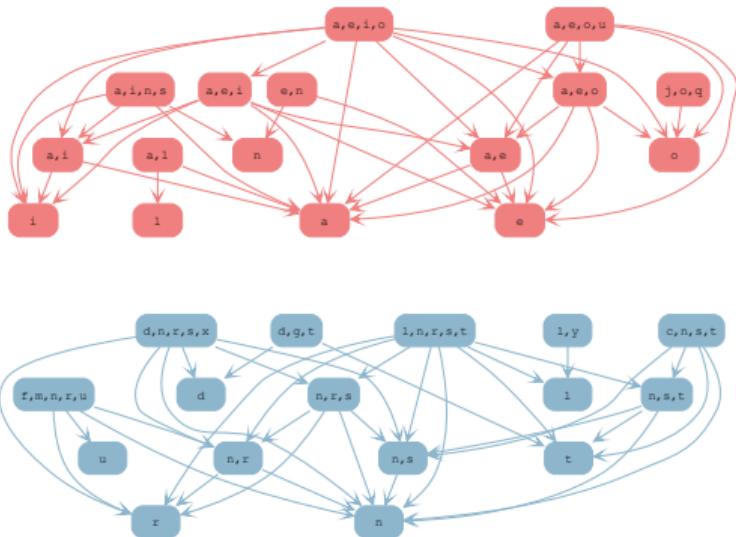
[G]iven a binary operation, noted

$a, b \rightsquigarrow \langle a|b \rangle : A \times B \rightarrow C$ and a subset $P \subset C$ (the ‘pole’) one can define the *polar* $X^\perp \subset B$ of a subset $X \subset A$ (resp. $Y^\perp \subset A$ of a subset $Y \subset B$) by :

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

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

- ◊ The map ‘polar’ is decreasing:
 $X \subset X' \Rightarrow X'^\perp \subset X^\perp.$
 - ◊ The set $\text{Pol}(A) \subset \mathcal{P}(A)$ of polar sets, i.e., of the form Y^\perp , is closed under arbitrary intersections. In particular, A is polar and $X^{\perp\perp}$ is the smallest polar set containing X .
 - ◊ As a consequence, $X^{\perp\perp\perp} = X^\perp$.



Theory of Computational Types

Definition (Polar/Orthogonal - Girard, 2011)

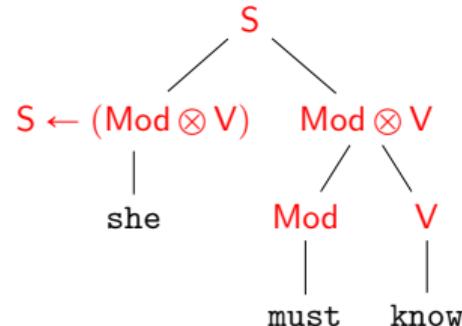
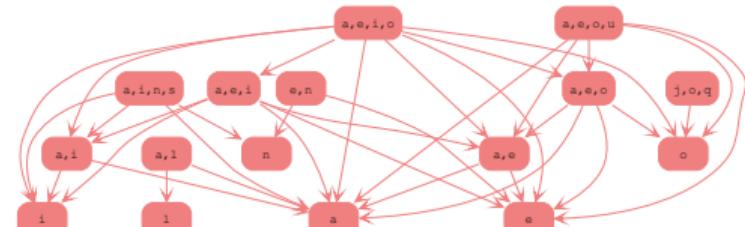
[G]iven a binary operation, noted

$a, b \rightsquigarrow \langle a|b \rangle : A \times B \rightarrow C$ and a subset $P \subset C$ (the 'pole')
 one can define the *polar* $X^\perp \subset B$ of a subset $X \subset A$
 (resp. $Y^\perp \subset A$ of a subset $Y \subset B$) by :

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

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

- ◊ The map 'polar' is decreasing:
 $X \subset X' \Rightarrow X'^\perp \subset X^\perp$.
- ◊ The set $\text{Pol}(A) \subset \mathcal{P}(A)$ of *polar* sets, i.e., of the form Y^\perp , is closed under arbitrary intersections. In particular, A is polar and $X^{\perp\perp}$ is the smallest polar set containing X .
- ◊ As a consequence, $X^{\perp\perp\perp} = X^\perp$.



(Gastaldi and Pellissier, 2021)

Outline

Introduction

NLMs as Formal Objects

The Structure(s) of the Embeddings

 The Algebra Behind the Embeddings

 The Structure Behind the Algebra

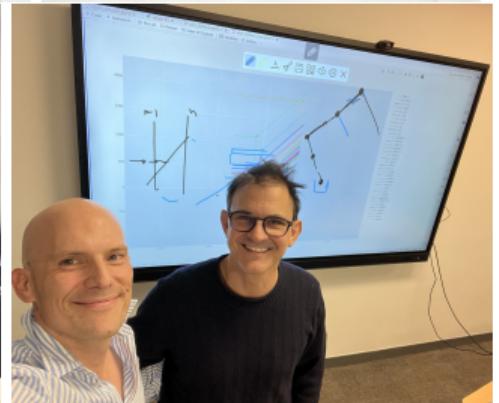
 The Categories Behind the Structure

Conclusion

Conclusion: For a Critical Formalism

- ◊ It is urgent to address of 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)

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*. <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>

References |

- Belrose, N., Schneider-Joseph, D., Ravfogel, S., Cotterell, R., Raff, E., & Biderman, S. (2024). Leace: Perfect linear concept erasure in closed form. *Proceedings of the 37th International Conference on Neural Information Processing Systems*.
- 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>
- 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.
- 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*. <https://doi.org/10.1080/03080188.2021.1890484>
- Girard, J.-Y. (2011, September). *The blind spot*. European Mathematical Society.

References II

- Gödel, K. (1986 (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.
- Kirschenbaum, M. (2023). *Again theory: A forum on language, meaning, and intent in the time of stochastic parrot*. <https://critiq.wordpress.com/2023/06/26/again-theory-a-forum-on-language-meaning-and-intent-in-the-time-of-stochastic-parrots/>
- 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., Le, Q., & Strohmann, T. (2013). *Learning representations of text using neural networks. NIPS deep learning workshop 2013 slides*.
- 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.
- 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://critiq.wordpress.com/2023/06/29/the-empirical-triumph-of-theory/>

Vernacular AI
Digital Theory Lab, NYU
New York, NY, USA

The Structure, Not the Prompt
For a Critical Formalism

Juan Luis Gastaldi

ETH zürich

February 7, 2025