



In Search of the Computational Primitives of Language

Aniello De Santo

he/him

MeLo Lab

Dept. of Linguistics

`aniellodesanto.github.io`

`aniello.desanto@utah.edu`

Computation and Theory Building

*[...] this is a confusion of two quite separate issues, **simulation and explanation**. [...] What we are **really** interested in [...] is explanation — in developing models that help us **understand how it is that people behave** that way, not merely demonstrating that we can build an artifact that behaves similarly.*

(Kaplan, 1995)

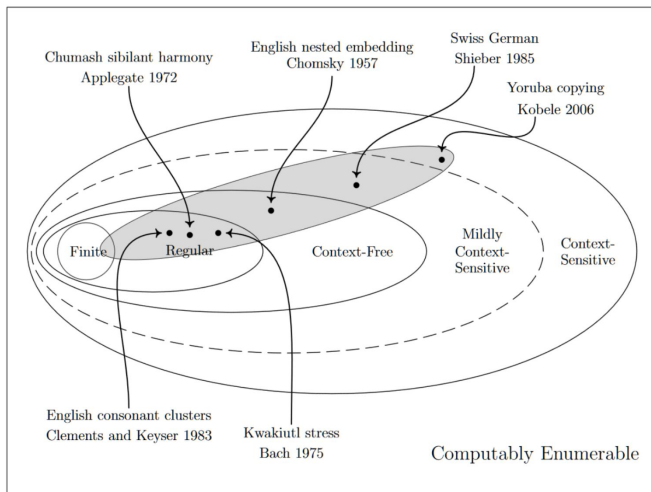
Computation and Theory Building

*[...] this is a confusion of two quite separate issues, **simulation and explanation**. [...] What we are **really** interested in [...] is explanation — in developing models that help us **understand how it is that people behave** that way, not merely demonstrating that we can build an artifact that behaves similarly.*

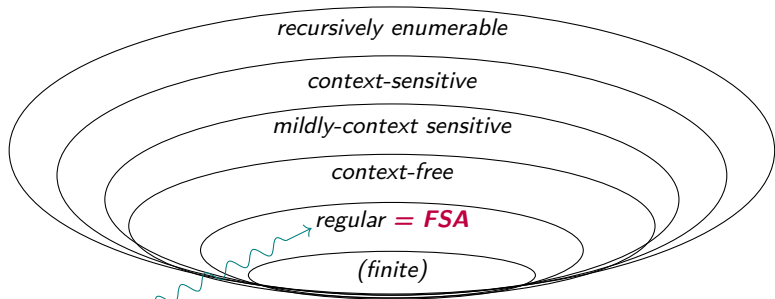
(Kaplan, 1995)

Interpretability for the win!

A Lens: Formal Language Theory



Spoken Languages' Phonology as a Regular System



Phonology

Kaplan and Kay (1994)

Local Phonotactic Dependencies

1 Intervocalic voicing in Italian

Forbid voiceless segments in between two vowels

- (1) a. * /kasa/
b. /kaza/
→ cf. orthography: “*casa*”

Intervocalic voicing is Strictly Local (SL)

- ▶ Forbid voiceless segments in-between two vowels: *V[-voice]V
- ▶ **Italian:** *ase, *ise, *ese, *isi, ...

\$ k a s a \$

\$ k a z a \$

Local Phonotactic Dependencies

1 Intervocalic voicing in Italian

Forbid voiceless segments in between two vowels

- (1) a. * /kasa/
b. /kaza/
→ cf. orthography: “*casa*”

Intervocalic voicing is Strictly Local (SL)

- ▶ Forbid voiceless segments in-between two vowels: *V[-voice]V
- ▶ **Italian:** *a**s**e, *i**s**e, *e**s**e, *i**s**i, ...

\$ k a **s** a \$

\$ k a z a \$

Local Phonotactic Dependencies

1 Intervocalic voicing in Italian

Forbid voiceless segments in between two vowels

- (1) a. * /kasa/
b. /kaza/
→ cf. orthography: “*casa*”

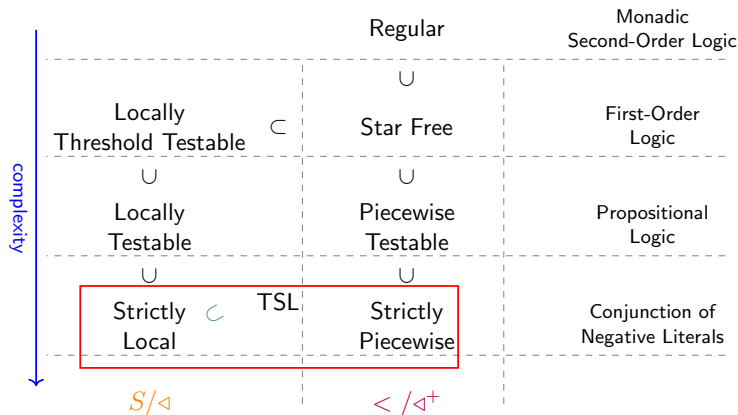
Intervocalic voicing is Strictly Local (SL)

- ▶ Forbid voiceless segments in-between two vowels: *V[-voice]V
- ▶ **Italian:** *a**s**e, *i**s**e, *e**s**e, *i**s**i, ...

* \$ k a **s** a \$

ok \$ k a z a \$

Beyond Automata: Subregular Languages¹



¹McNaughton & Papert (1976), Heinz (2011), Chandlee & Heinz (2014), De Santo & Graf (2019), De Santo & Rawski (2022), a.o.

FLT, Linguistics, and LLM Expressivity

Saturated Transformers are Constant-Depth Threshold Circuits

William Merrill^{*†} Ashish Sabharwal^{*} Noah A. Smith^{*‡}
^{*} Allen Institute for AI [†] New York University [‡] University of Washington
{ashishs,noah}@allenai.org

Between Circuits and Chomsky: Pre-pretraining on Formal Languages Imparts Linguistic Biases

Michael Y. Hu¹ Jackson Petty² Chuan Shi¹ William Merrill¹ Tal Linzen^{1,2}

Neuro-Symbolic Language Modeling with Automaton-augmented Retrieval

Uri Alon¹ Frank F. Xu¹ Junxian He¹ Sudipta Sengupta² Dan Roth³ Graham Neubig¹
¹Language Technologies Institute, Carnegie Mellon University ²Amazon AWS ³AWS AI Labs
{ualon,fangzhex,junxianh,gneubig}@cs.cmu.edu {sudipta,drot}@amazon.com

What Makes Instruction Learning Hard? An Investigation and a New Challenge in a Synthetic Environment

Matthew Finlayson Kyle Richardson Ashish Sabharwal Peter Clark
Allen Institute for AI, Seattle, WA
{matthewf,kyler,ashishs,peterc}@allenai.org



Switching Gears: Benchmarking with Sentence Processing

Not All Structures Are Processed Equally

- Subject VS object relative clause

SRC **The horse** [_{RC} that kicked the wolf] went home.

ORC **The horse** [_{RC} that the wolf kicked] went home.

Switching Gears: Benchmarking with Sentence Processing

Not All Structures Are Processed Equally

- Subject VS object relative clause

SRC **The horse** [_{RC} that kicked the wolf] went home.

ORC **The horse** [_{RC} that the wolf kicked] went home.

Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies

Tal Linzen^{1,2} **Emmanuel Dupoux**¹
LSCP¹ & IJN², CNRS,
EHESS and ENS, PSL Research University
{tal.linzen,
emmanuel.dupoux}@ens.fr

Yoav Goldberg
Computer Science Department
Bar Ilan University
yoav.goldberg@gmail.com

Switching Gears: Benchmarking with Sentence Processing

Not All Structures Are Processed Equally

- Subject VS object relative clause

SRC **The horse** [_{RC} that **t** kicked the wolf] went home.

ORC **The horse** [_{RC} that the wolf kicked **t**] went home.

- Attachment preferences

1.a I shot an elephant in my pajamas

1.b I [shot an elephant] [in my pajamas]

Switching Gears: Benchmarking with Sentence Processing

Not All Structures Are Processed Equally

- ▶ Subject VS object relative clause

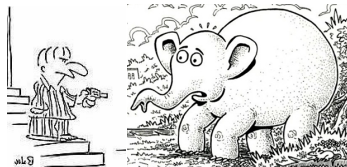
SRC **The horse** [_{RC} that **t** kicked the wolf] went home.

ORC **The horse** [_{RC} that the wolf kicked **t**] went home.

- ▶ Attachment preferences

1.a I shot **[an elephant in my pajamas]**

1.b I **[shot an elephant]** **[in my pajamas]**



Ambiguity All the Way Down



So What?

Ambiguity is ubiquitous in natural language!

For Cognitive Science

- ▶ How do humans handle multiple structural representations?
- ▶ What principles guide ambiguity resolution cross-linguistically?
- ▶ Language specific properties vs. general biases/mechanisms?

For NLP

- ▶ How do LLMs handle multiple structural representations?
- ▶ What principles guide ambiguity resolution cross-linguistically?
- ▶ Language specific properties vs. general biases/mechanisms?

So What?

Ambiguity is ubiquitous in natural language!

For Cognitive Science

- ▶ How do humans handle multiple structural representations?
- ▶ What principles guide ambiguity resolution cross-linguistically?
- ▶ Language specific properties vs. general biases/mechanisms?

For NLP

- ▶ How do LLMs handle multiple structural representations?
- ▶ What principles guide ambiguity resolution cross-linguistically?
- ▶ Language specific properties vs. general biases/mechanisms?

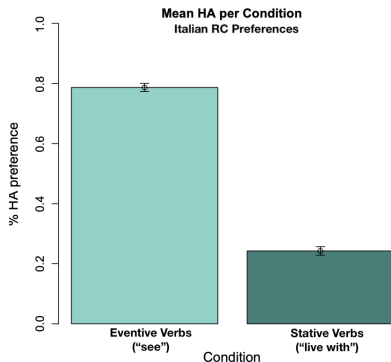
Ambiguity and Relative Clauses Cross-linguistically ²

- ▶ They saw the daughter of the actress that was on the balcony
- | | | | |
|-----------------|---------------------|--------------------|----|
| NP ₁ | The daughter | was on the balcony | HA |
| NP ₂ | The actress | was on the balcony | LA |

²Grillo & Costa (2015,) De Santo & Lee (2023), Lee & De Santo (2024)

Ambiguity and Relative Clauses Cross-linguistically ²

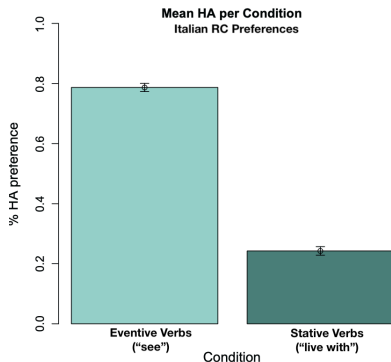
- ▶ They saw the daughter of the actress that was on the balcony
- | | | | |
|-----------------|---------------------|--------------------|----|
| NP ₁ | The daughter | was on the balcony | HA |
| NP ₂ | The actress | was on the balcony | LA |



²Grillo & Costa (2015,) De Santo & Lee (2023), Lee & De Santo (2024)

Ambiguity and Relative Clauses Cross-linguistically ²

- ▶ They saw the daughter of the actress that was on the balcony
- | | | | |
|-----------------|---------------------|--------------------|----|
| NP ₁ | The daughter | was on the balcony | HA |
| NP ₂ | The actress | was on the balcony | LA |

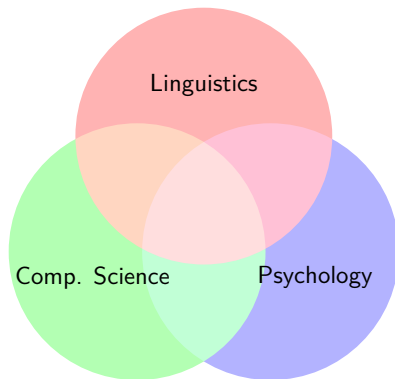


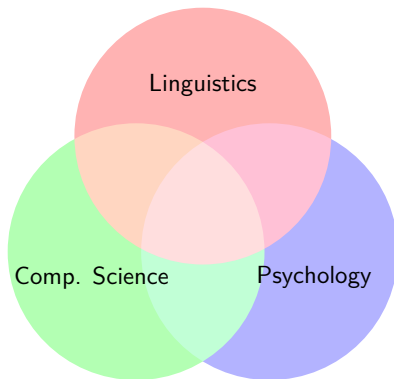
LLMs and Italian RCs?



Check out the poster!

²Grillo & Costa (2015,) De Santo & Lee (2023), Lee & De Santo (2024)





Let's chat!

Unbounded Dependencies Are Not SL

► **Samala Sibilant Harmony**

Sibilants must not disagree in anteriority.

(?)

- (2) a. * ha **s**xintilawa **ʃ**
- b. * ha **ʃ**xintilawa **s**
- c. ha **ʃ**xintilawa **ʃ**

Example: Samala

*\$ h a **s** x i n t i l a w a **ʃ** \$

\$ h a **ʃ** x i n t i l a w a **ʃ** \$

Unbounded Dependencies Are Not SL

► Samala Sibilant Harmony

Sibilants must not disagree in anteriority.

(?)

- (2) a. * ha **s** x i n t i l a w a **ʃ**
- b. * ha **ʃ** x i n t i l a w a **s**
- c. ha **ʃ** x i n t i l a w a **ʃ**

Example: Samala

* \$ h a **s** x i n t i l a w a **ʃ** \$

\$ h a **ʃ** x i n t i l a w a **ʃ** \$

Unbounded Dependencies Are Not SL

► Samala Sibilant Harmony

Sibilants must not disagree in anteriority.

(?)

- (2) a. * ha^sxintilawa^f
b. * ha^fxintilawa^s
c. ha^fxintilawa^f

Example: Samala

* \$ h a ^s x i n t i l a w a ^f \$

\$ h a ^f x i n t i l a w a ^f \$

Unbounded Dependencies Are Not SL

► Samala Sibilant Harmony

Sibilants must not disagree in anteriority.

(?)

- (2) a. * ha^sxintilawa^f
b. * ha^fxintilawa^s
c. ha^fxintilawa^f

Example: Samala

*\$ ha^sxintilawa^f\$
\$ ha^fxintilawa^f\$

Unbounded Dependencies Are Not SL

► Samala Sibilant Harmony

Sibilants must not disagree in anteriority.

(?)

- (2) a. * ha^sxintilawa^ʃ
b. * ha^ʃxintilawa^s
c. ha^ʃxintilawa^ʃ

Example: Samala

* \$ ha^s x i n t i l a w a ʃ \$
\$ ha^ʃ x i n t i l a w a ʃ \$

► **But:** Sibilants can be arbitrarily far away from each other!

* \$ ^s t a j a n o w o n w a ʃ \$

Unbounded Dependencies Are Not SL

► Samala Sibilant Harmony

Sibilants must not disagree in anteriority.

(?)

- (2) a. * ha^sxintilawa^f
b. * ha^fxintilawa^s
c. ha^fxintilawa^f

Example: Samala

* \$ ha^sxintilawa^f \$
\$ ha^fxintilawa^f \$

► **But:** Sibilants can be arbitrarily far away from each other!

* \$ ^sstajanowonwa^f \$

Unbounded Dependencies are TSL

- ▶ Let's revisit Samala Sibilant Harmony

- (3) a. * ha^sxintilawa^ʃ
b. * ha^ʃxintilawa^s
c. ha^ʃxintilawa^ʃ

- ▶ What do we need to project? [+strident]
- ▶ What do we need to ban? *[+ant][−ant], *[−ant][+ant]

I.E. *^sʃ, *^sʒ, *^ʃʃ, *^ʃʒ, *^ʃs, *^ʒs, *^ʃz, *^ʒz

Example: TSL Samala

^s	ʃ		ʃ	ʃ
.....			
* \$ha ^s xintilawa ^ʃ \$			<i>ok</i> \$ha ^ʃ xintilawa ^ʃ \$	

Unbounded Dependencies are TSL


- ▶ Let's revisit Samala Sibilant Harmony

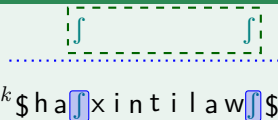
- (3) a. * ha^sxintilawa^ʃ
b. * ha^ʃxintilawa^s
c. ha^ʃxintilawa^ʃ

- ▶ What do we need to project? [+strident]
- ▶ What do we need to ban? *[+ant][−ant], *[−ant][+ant]

I.E. *^sʃ, *^sʒ, *^ʒʃ, *^ʒʒ, *^ʃs, *^ʒs, *^ʃz, *^ʒz

Example: TSL Samala


* \$ha^sxintilawa^ʃ\$


ok \$ha^ʃxintilawa^ʃ\$