

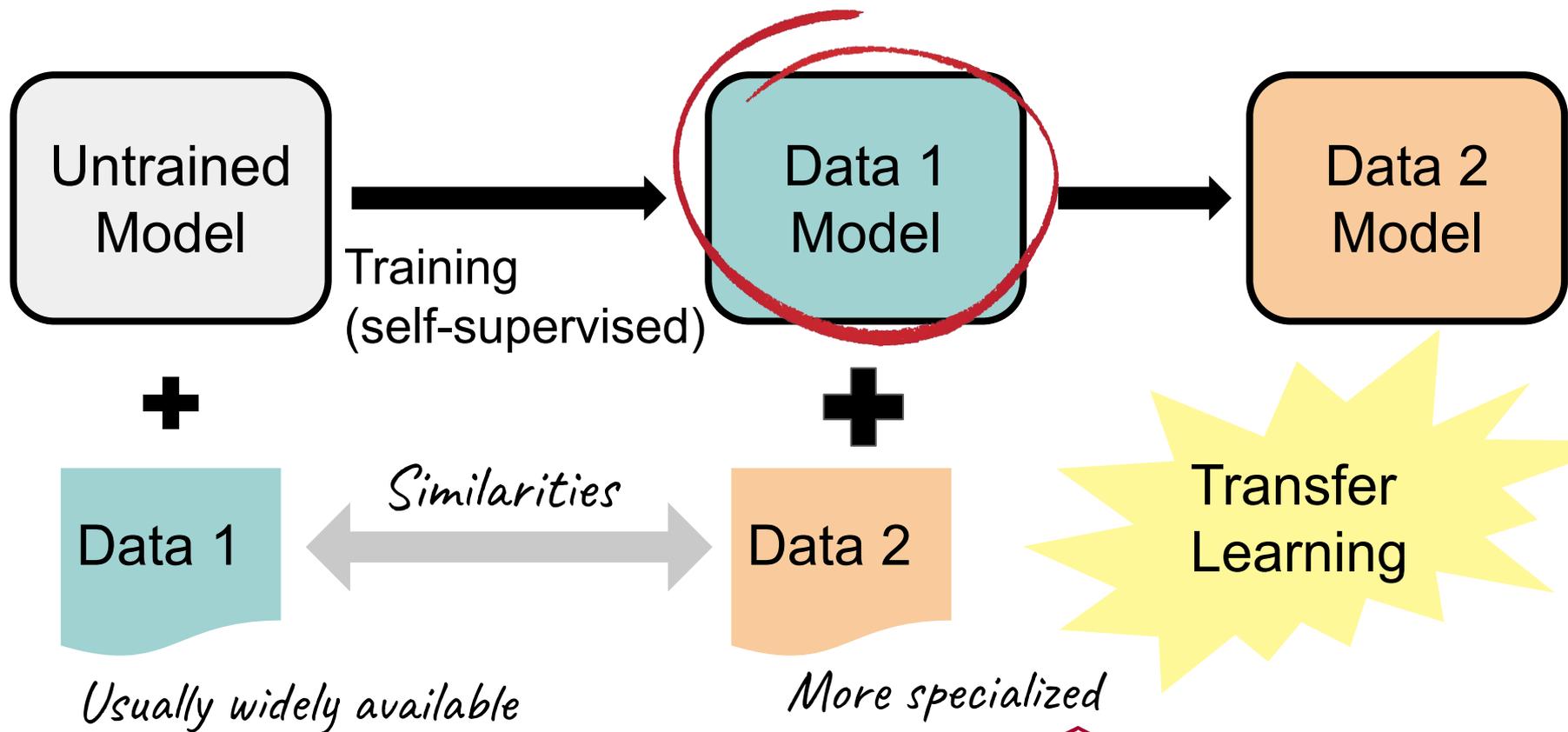
Structural Transfer Learning: Exploring Neural Models and Language Structure Through Understanding Transfer



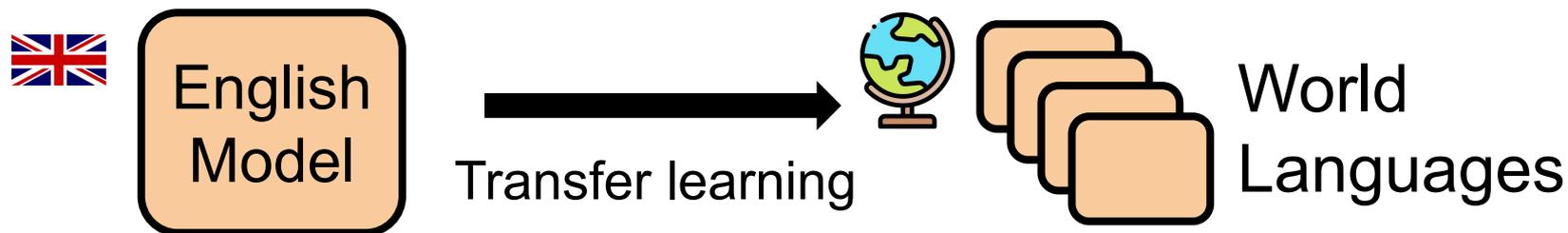
Isabel Papadimitriou



Transfer Learning



Transfer learning has practical applications



But also an **analysis methodology** for understanding data and learning

- Power machine learning models let us explore questions about language in new ways

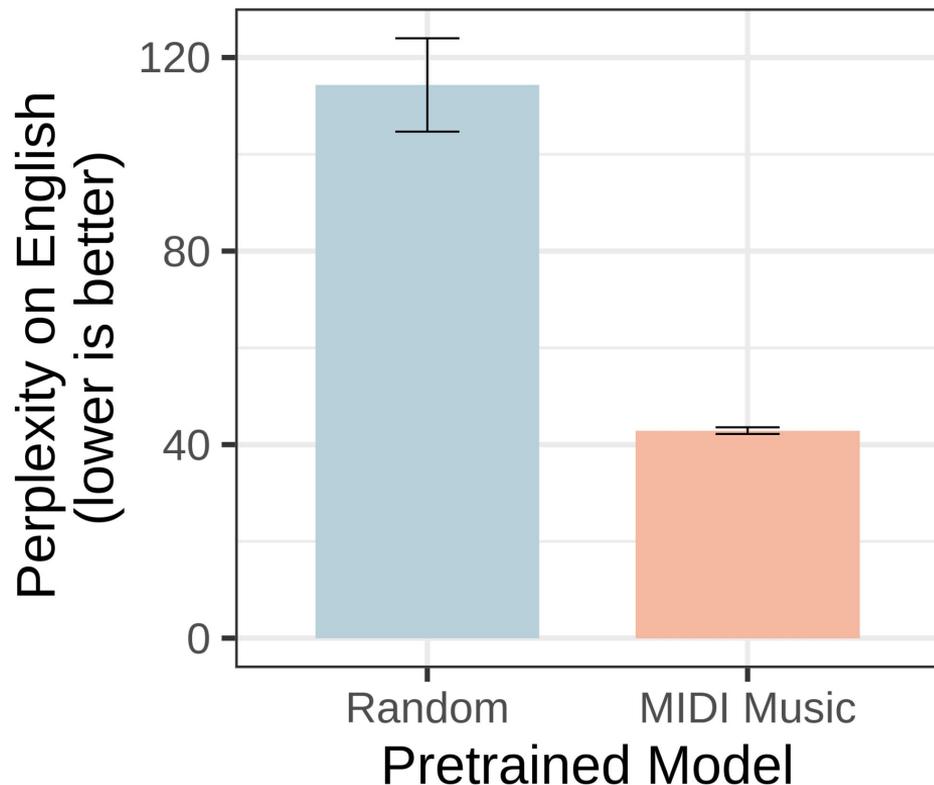
Shared structures between modalities

Data 1

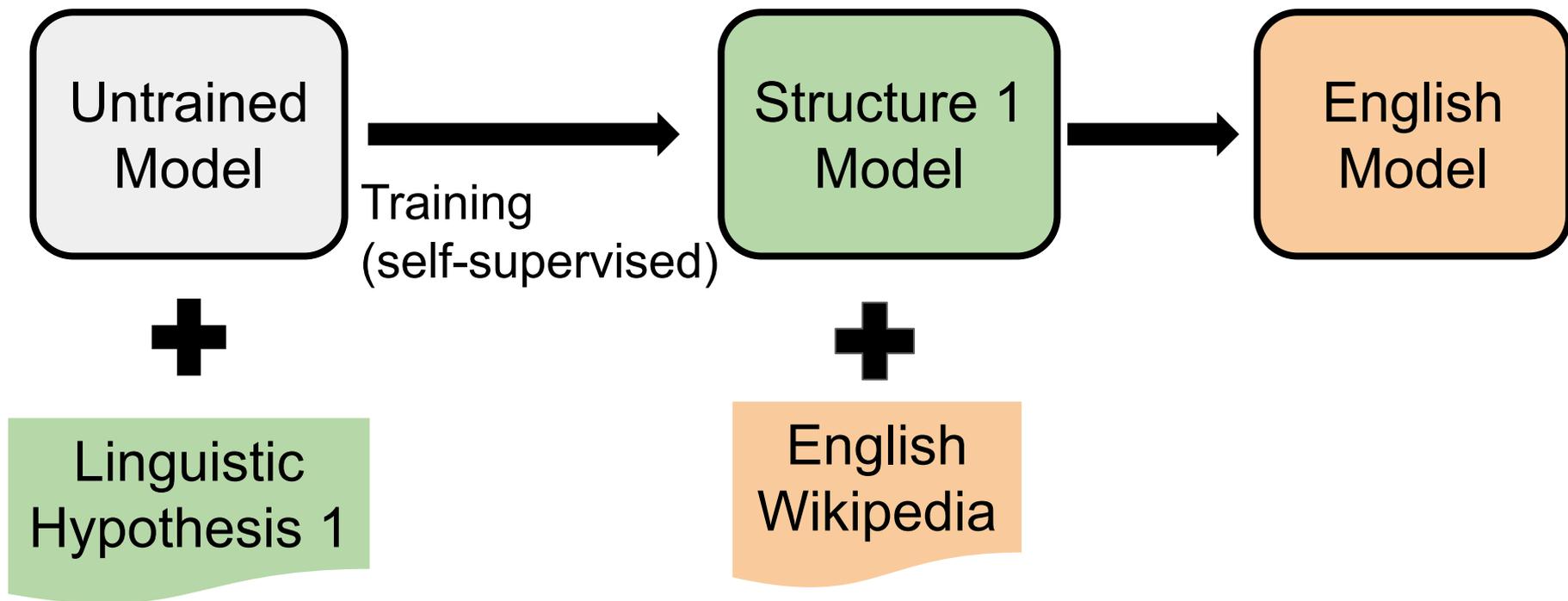
Random numbers
or
MIDI Music

Data 2

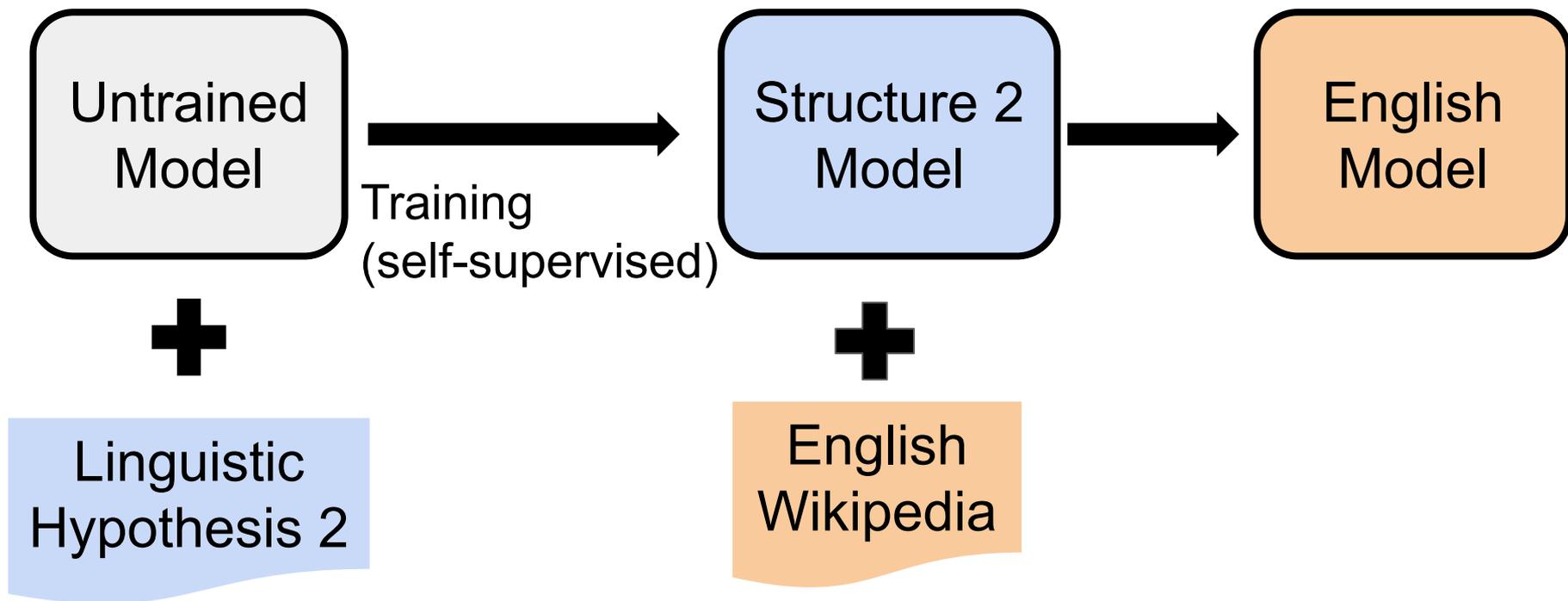
English wikipedia
(~100 mil token
subset)



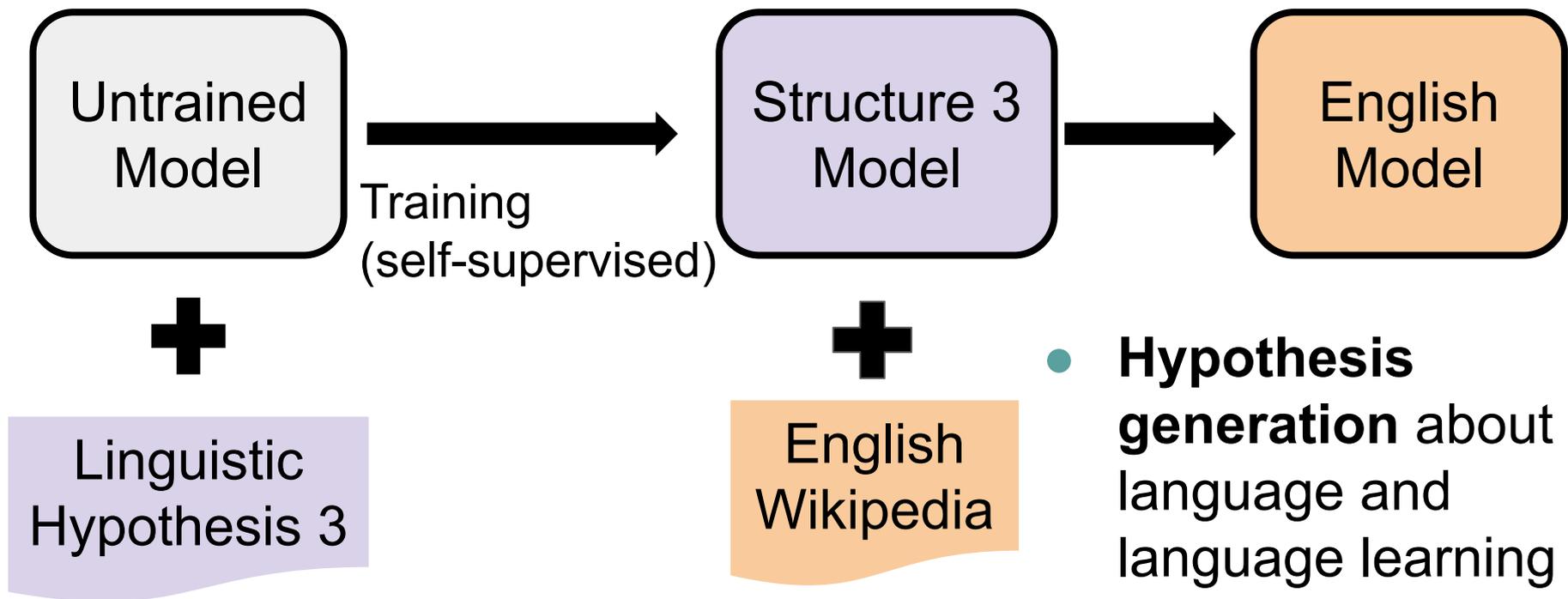
Structural Transfer: a testbed for linguistic structure hypotheses

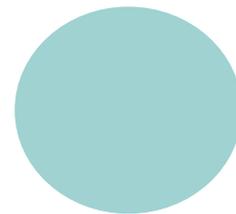


Structural Transfer: a testbed for linguistic structure hypotheses



Structural Transfer: a testbed for linguistic structure hypotheses





Using structural transfer learning to explore the **role of structure** in language and language learning

Transfer learning in NLP

- Recent NLP: pretrain so much, that the task can be described in language. **Prompting**

Transfer learning now

- Looking beyond the **dominant languages** where we can do things like prompting
- And for understanding **structure**

Structure and language

- Structure is characteristic of human language
- Most obviously in syntax
- But also beyond syntax
 - Meaning, discourse, reference, information structure
- What structural biases are sufficient for language learning?
- (beyond this talk) Role of communication and language use in creating structure

Three hypotheses about language

1) Recursion

Constituents “Clumping”

The cat sat on the mat

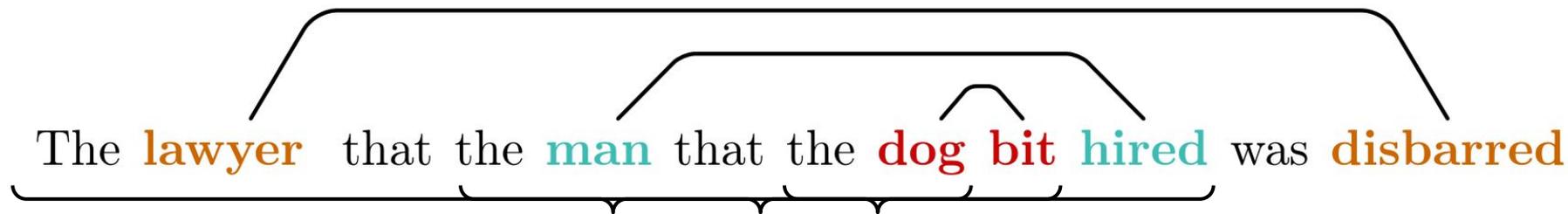
I think that the cat sat on the mat

You always accuse me that I think that the cat sat on the mat

Three hypotheses about language

1) Recursion

Nesting Context-free

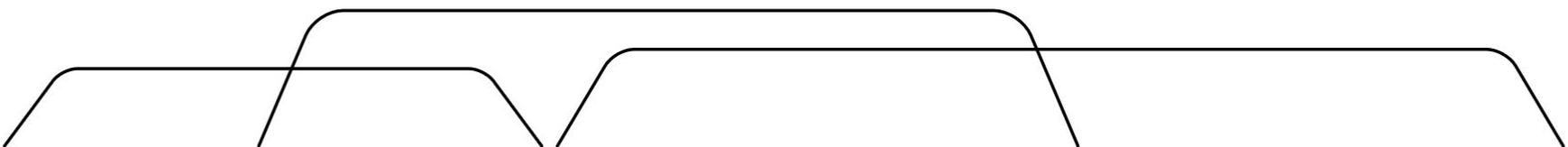


Three hypotheses about language

2) Crossing links and dependencies

Linking in meaning and reference

“I voted for **him** even though **I** am negatively affected by **his** redistribution policies” **he** said

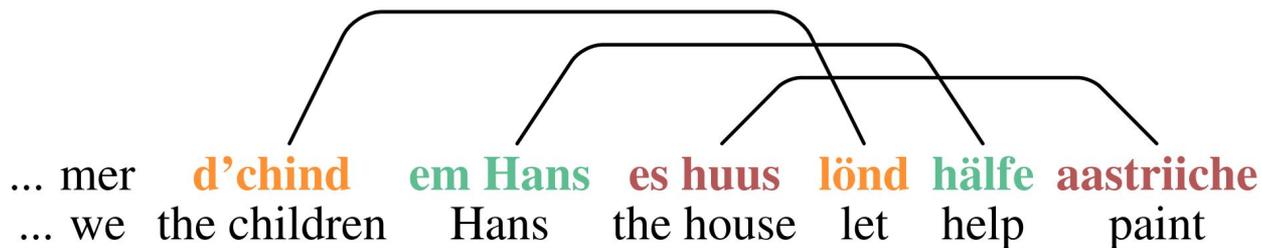


The diagram illustrates the crossing dependencies between pronouns in the sentence. It features three horizontal lines representing dependencies. The first line connects the first 'I' to 'he'. The second line connects 'him' to the second 'I'. The third line connects 'his' to the first 'I'. These lines cross each other, demonstrating the complexity of linking pronouns in natural language.

Three hypotheses about language

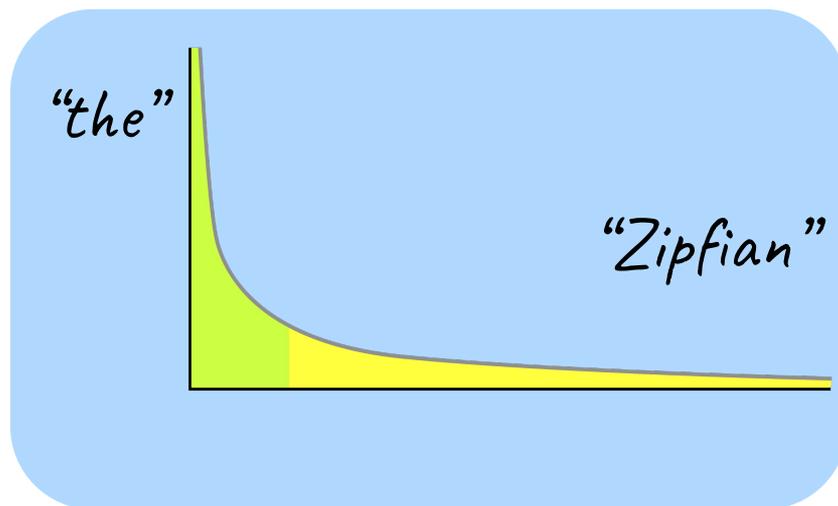
2) Crossing links and dependencies

And syntactic structures



Three hypotheses about language

3) Zipfian vocabulary distribution



Outline

- What **structural biases** are useful for human language learners?
 - Disentangling the effects of **recursive** and **linking** structures
- How does **vocabulary distribution** transfer as a structural bias?
 - The structural effect of vocabulary

Outline

- What **structural biases** are useful for human language learners?
 - Disentangling the effects of **recursive** and **linking** structures
- How does vocabulary distribution transfer as a structural bias?
 - The structural effect of vocabulary

Exploring inductive bias

?

+

Pretty limited
linguistic exposure



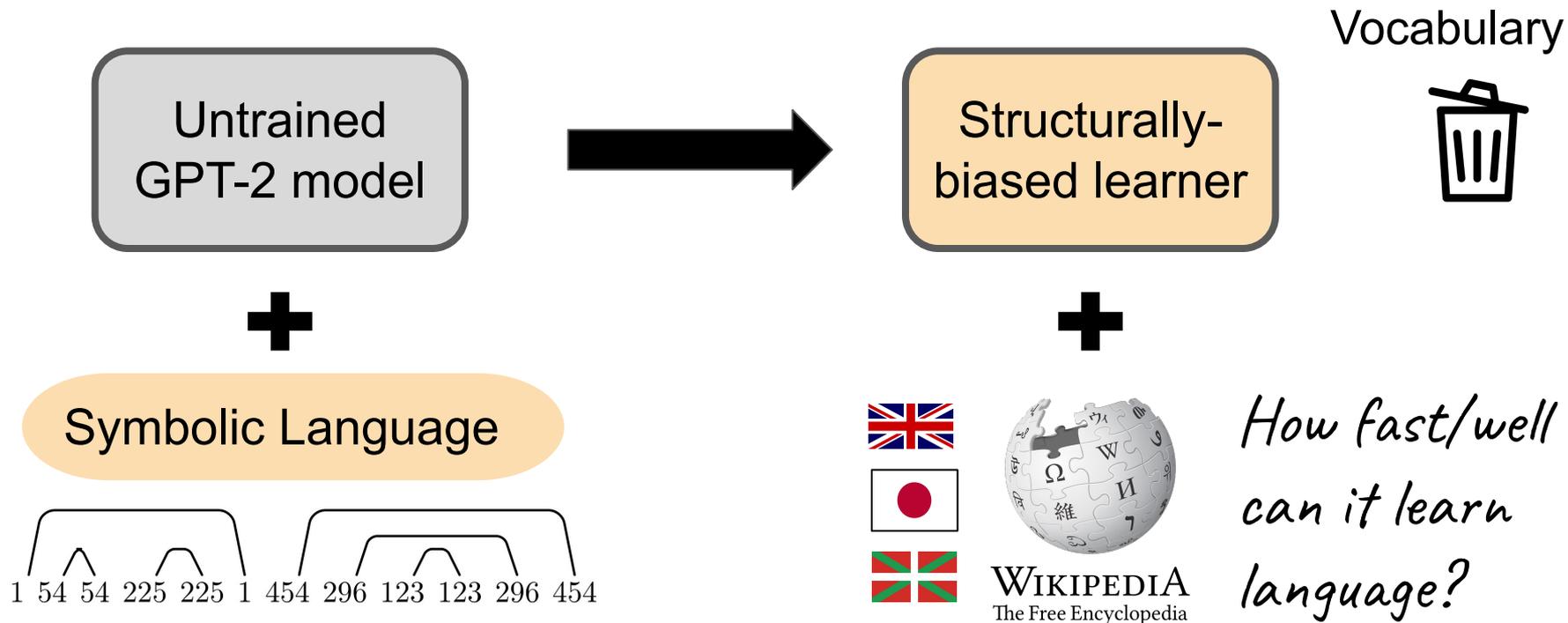
Language learning



Use transfer learning to test different **structural inductive learning biases**



Transfer learning methodology

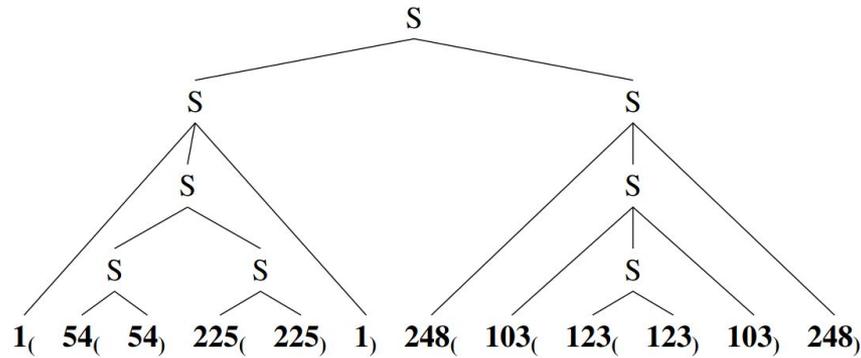


Symbolic pretraining languages

Nesting
Parentheses



Nested Parentheses Primitive



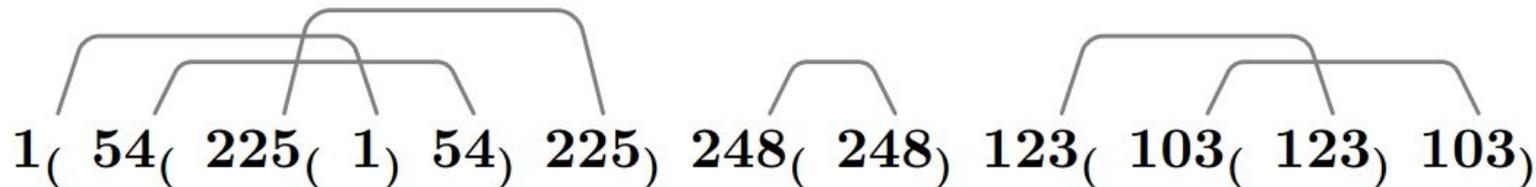
- Well-nested, matching pairs
- Constituents

Symbolic pretraining languages

Nesting
Parentheses

Crossing
Dependencies

Crossing Dependencies



- Tokens have to **match**, but not **nest**
- Where does the structure come from?
 - **Dependency length distribution:** sample from empirical distances of nesting parentheses

Symbolic pretraining languages

Nesting
Parentheses

Crossing
Dependencies

Random

Regular
repetition

Controls:



Simple Repetition Primitive

Randomly sample k words, then repeat them, then randomly sample k words...

499 472 300 345 272 **499 472 300 345 272** **309 17 15**

(Example is for $k=5$, we do $k=10$ in experiments)

Symbolic pretraining languages

Nesting
Parentheses

Crossing
Dependencies

Controls:

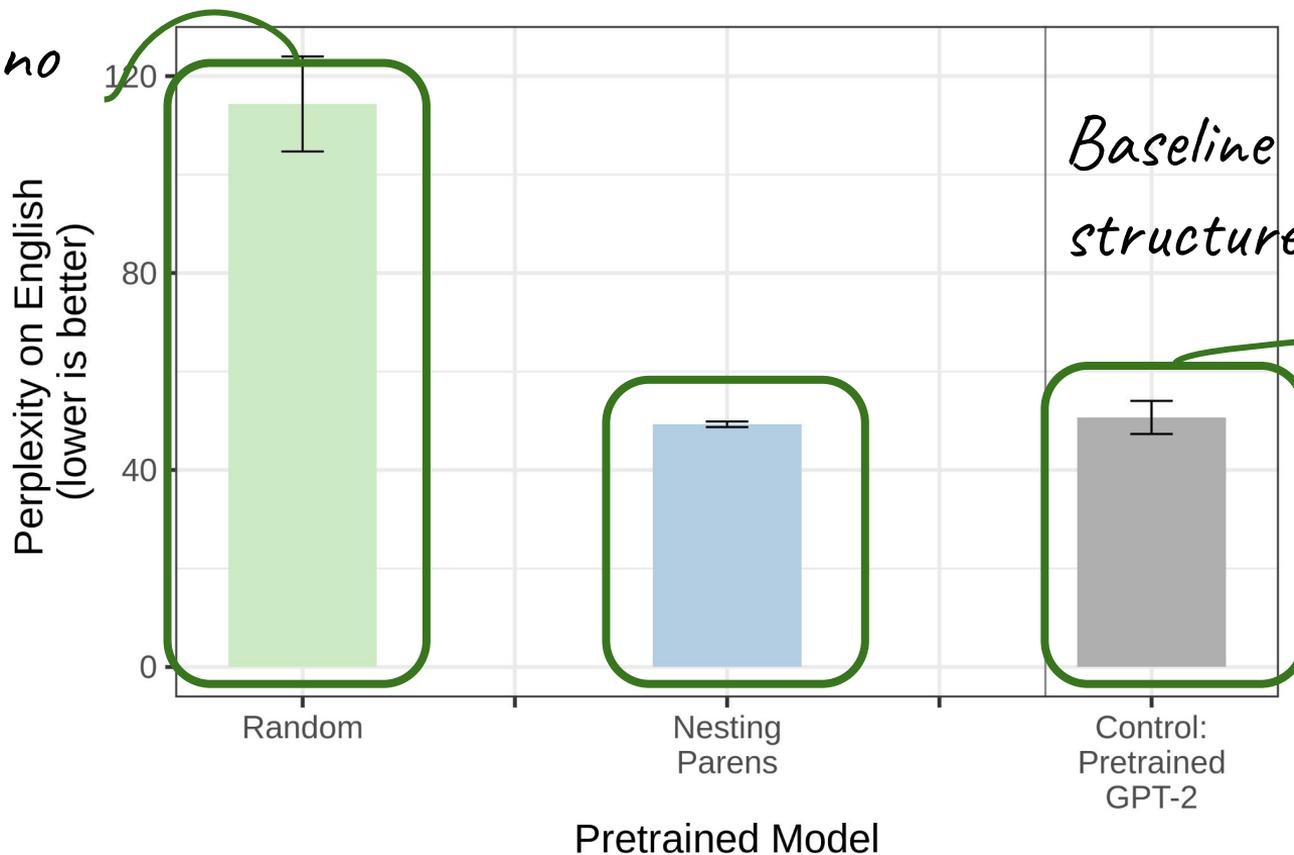
Random

Regular
repetition



Nesting structure helps language learning

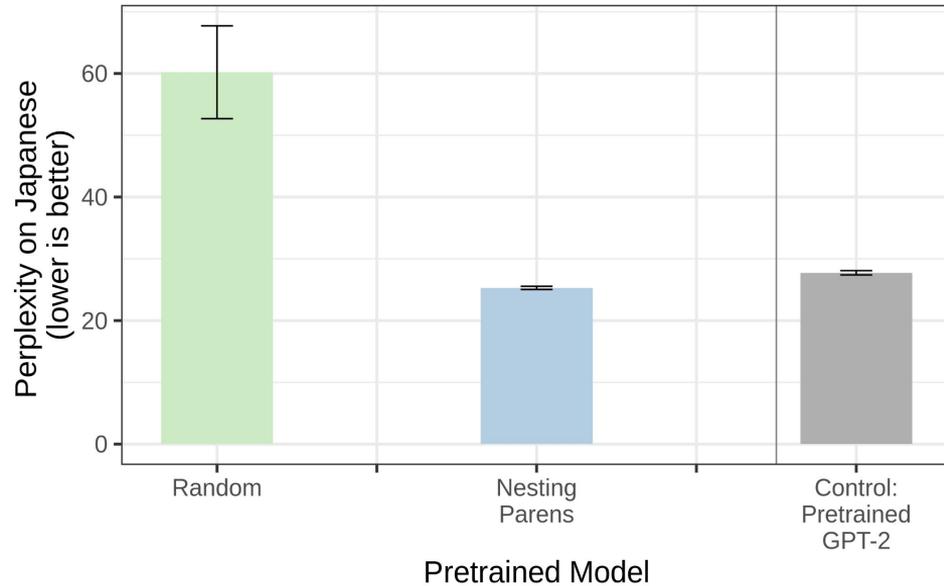
Baseline - no structure



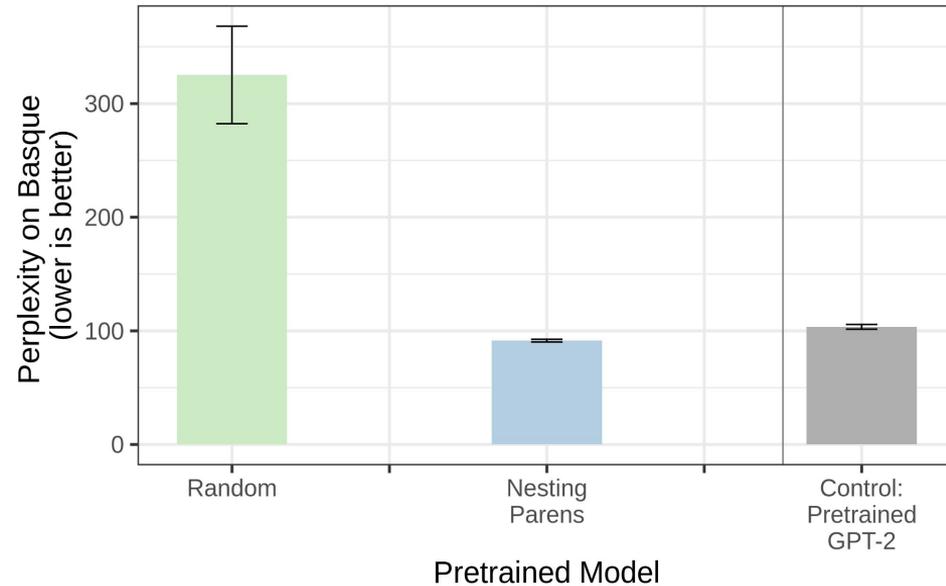
Baseline - English structure

Multilingual case – Japanese and Basque

Japanese



Basque



Symbolic pretraining languages

Nesting
Parentheses

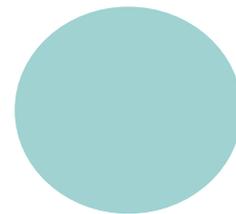
Crossing
Dependencies

Random

Regular
repetition

Controls:



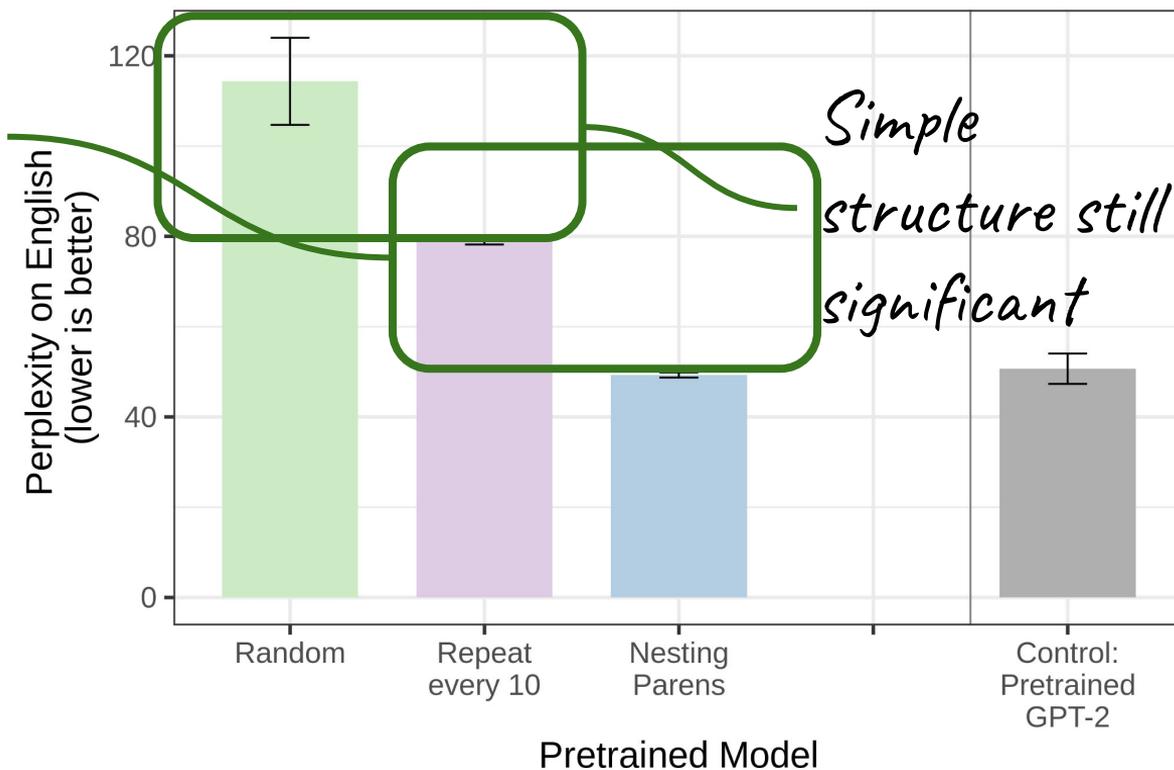


Question: does **nesting** really help? Or would any structure help?

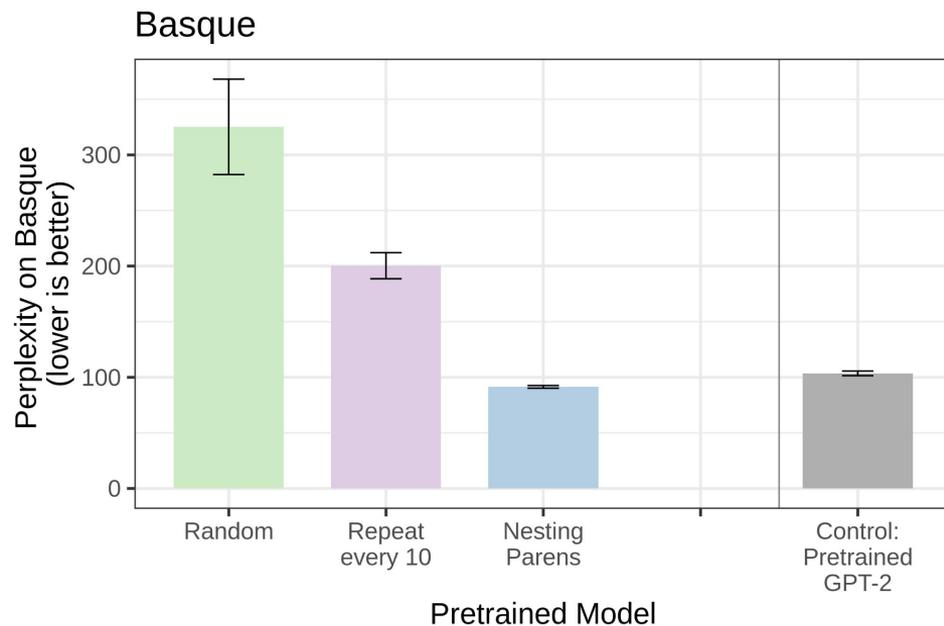
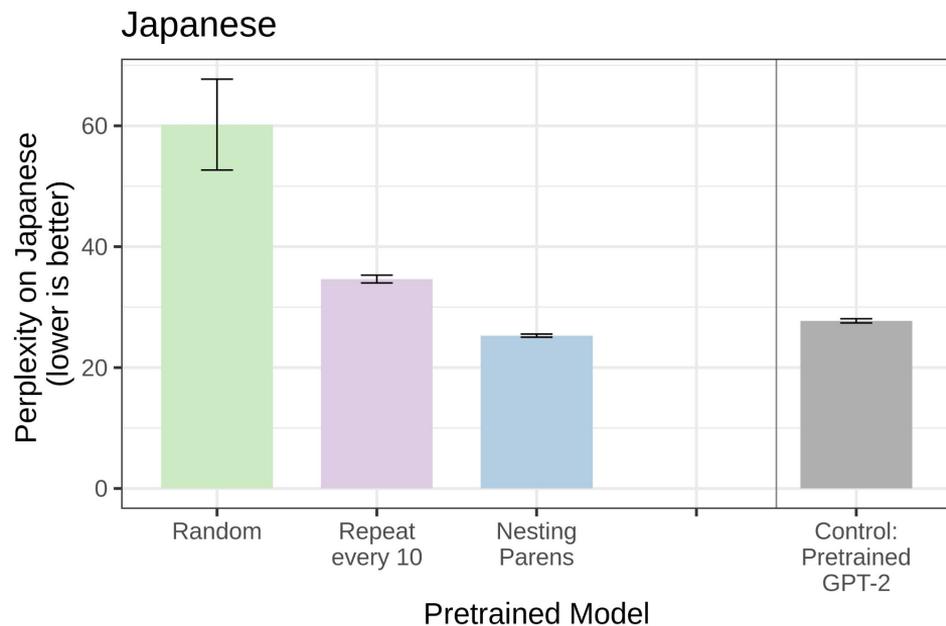
499 472 300 345 272 499 472 300 345 272 309 17 15

Not just any structure has this effect

Complexity of nesting is necessary



Again, a multilingual effect



Symbolic pretraining languages

Controls:

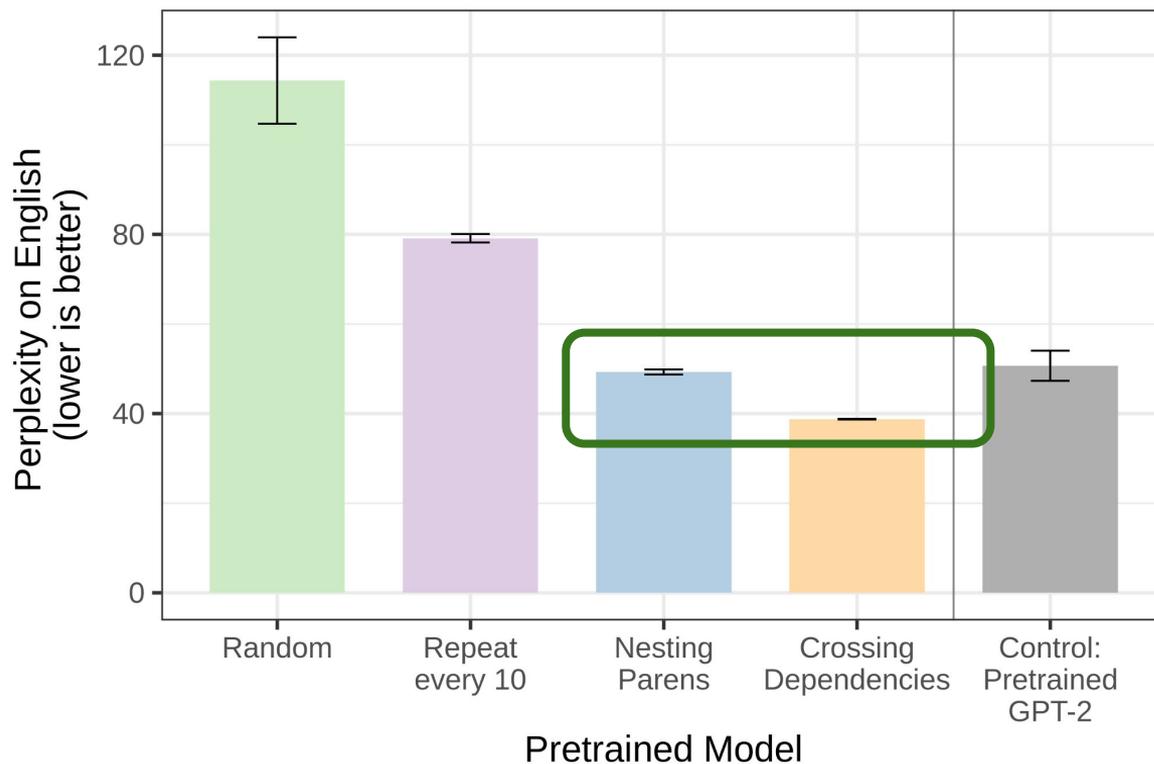
Nesting
Parentheses

Crossing
Dependencies

Random

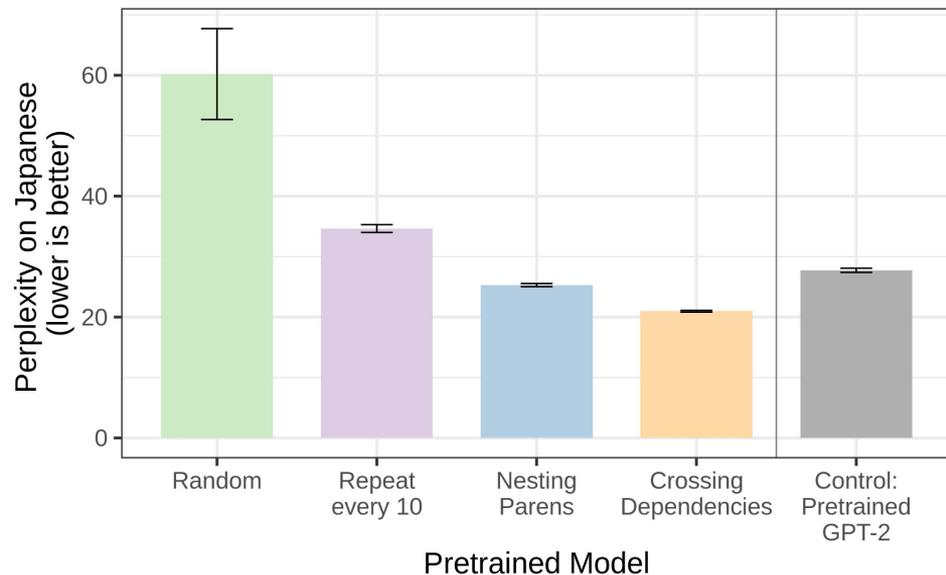
Regular
repetition

Crossing links, without nesting, provide a better inductive bias

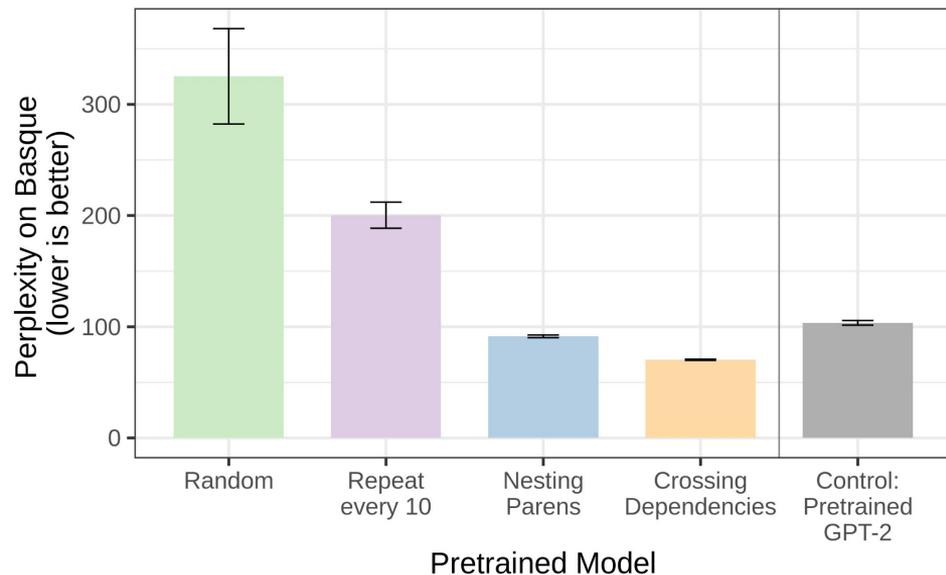


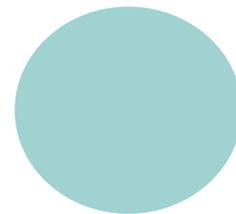
This is also true across languages

Japanese



Basque

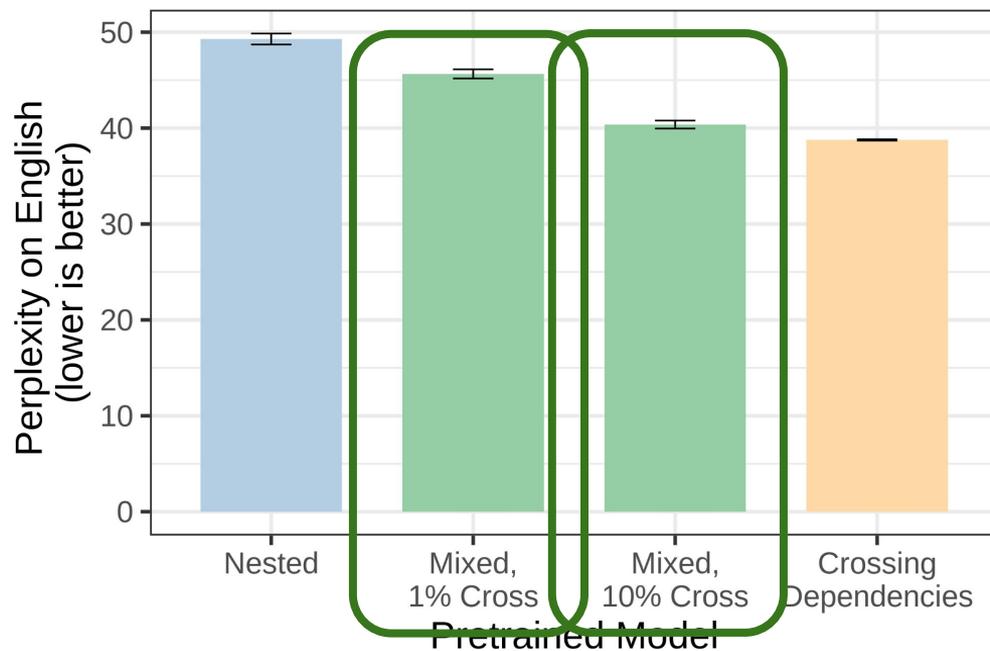




The kinds of structure that make language are multifaceted

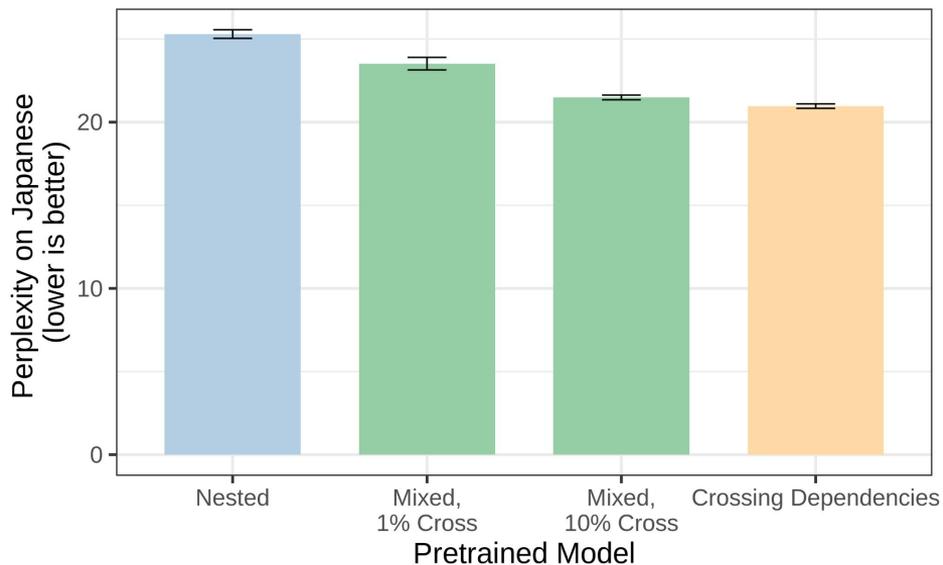
- Structural transfer lets us explore hypotheses about structure in language
- Language as a learnable system, independent of linguistic theory

Slightly breaking constituent structure makes better language learners

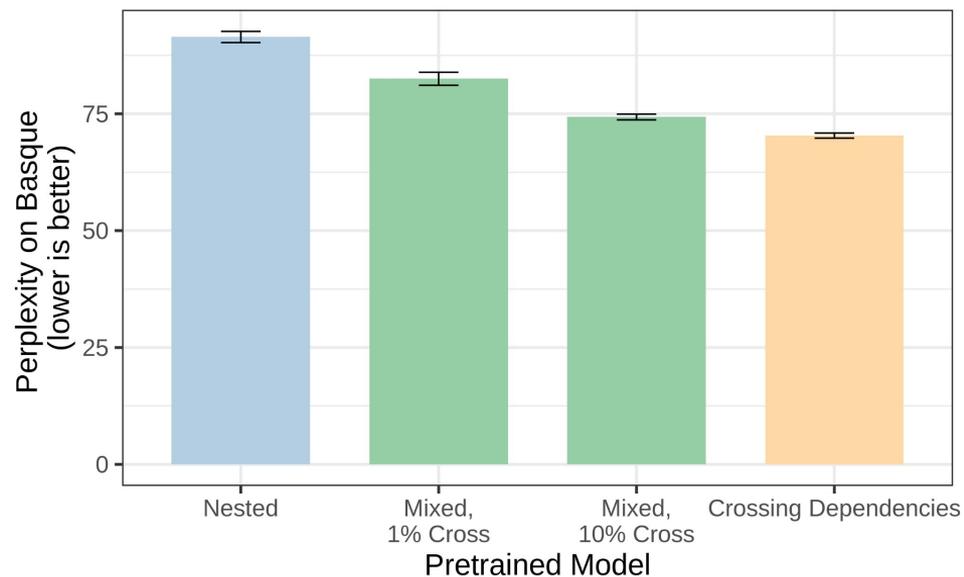


Also a multilingual effect

Japanese



Basque



Structural inductive bias through transfer learning

- Complex structural relationships are important in language
- Multiple crossing dependencies a good starting point for language learning
- Computational models as **hypothesis generators**: testing linguistic structure in theory-free ways

Outline

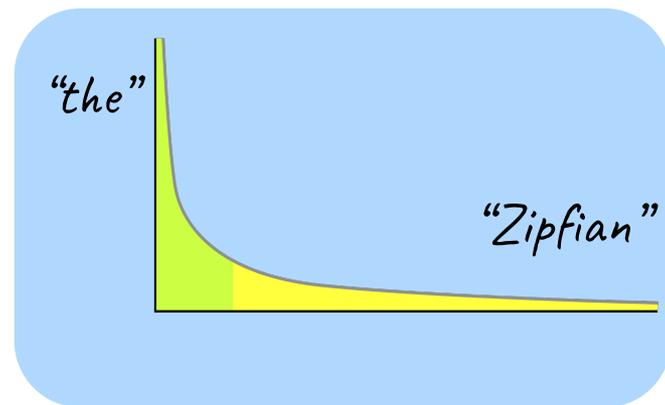
- What **structural biases** are useful for human language learners?
 - Disentangling the effects of **recursive** and **linking** structures
- How does vocabulary distribution transfer as a structural bias?
 - The structural effect of vocabulary

Outline

- What structural biases are useful for human language learners?
 - Disentangling the effects of **recursive** and **linking** structures
- How does **vocabulary distribution** transfer as a structural bias?
 - The structural effect of vocabulary

The lexicon in linguistics

- A good amount of structure is in the vocabulary:
- Vocabulary distribution
- Structure in meaning
- and also in grammar
 - Properties like transitive verb
 - Constructions, like “Let alone”
 - ...



We throw out the vocabulary between pretraining and fine-tuning

Pretraining vocabulary

(1 (2 (3 (4 (5 (6 (7 (8 (9 (10 (11 (12 (13 (14 (15 (16 (17 (18 (19 (20 (21 (22 (23 (24 (25 (26 (27 (28 (29 (30 (31 (32 (33 (34 (35 (36 (37 (38 (39 (40 (41 (42 (43 (44 (45 (46 (47 (48 (49 (50 (51 (52 (53 (54 (55 (56 (57 (58 (59 (60 (61 (62 (63 (64 (65 (66 (67 (68 (69 (70 (71 (72 (73 (74 (75 (76 (77 (78 (79 (80 (81 (82 (83 (84 (85 (86 (87 (88 (89 (90 (91 (92 (93 (94 (95 (96 (97 (98 (99 (100 (101 (102 (103 (104 (105 (106 (107 (108 (109 (110 (111 (112 (113 (114 (115 (116 (117 (118 (119 (120 (121 (122 (123 (124 (125 (126 (127 (128 (129 (130 (131 (132 (133 (134 (135 (136 (137 (138 (139 (140 (141 (142 (143 (144 (145 (146 (147 (148 (149 (150 (151 (152 (153 (154 (155 (156 (157 (158 (159 (160 (161 (162 (163 (164 (165 (166 (167 (168 (169 (170 (171 (172 (173 (174 (175 (176 (177 (178 (179 (180 (181 (182 (183 (184 (185 (186 (187 (188 (189 (190 (191 (192 (193 (194 (195 (196 (197 (198 (199 (200 (201 (202 (203 (204 (205 (206 (207 (208 (209 (210 (211 (212 (213 (214 (215 (216 (217 (218 (219 (220 (221 (222 (223 (224 (225 (226 (227 (228 (229 (230 (231 (232 (233 (234 (235 (236 (237 (238 (239 (240 (241 (242 (243 (244 (245 (246 (247 (248 (249 (250

Vocabulary indices 0-499

What the model sees

Fine-tuning vocabulary

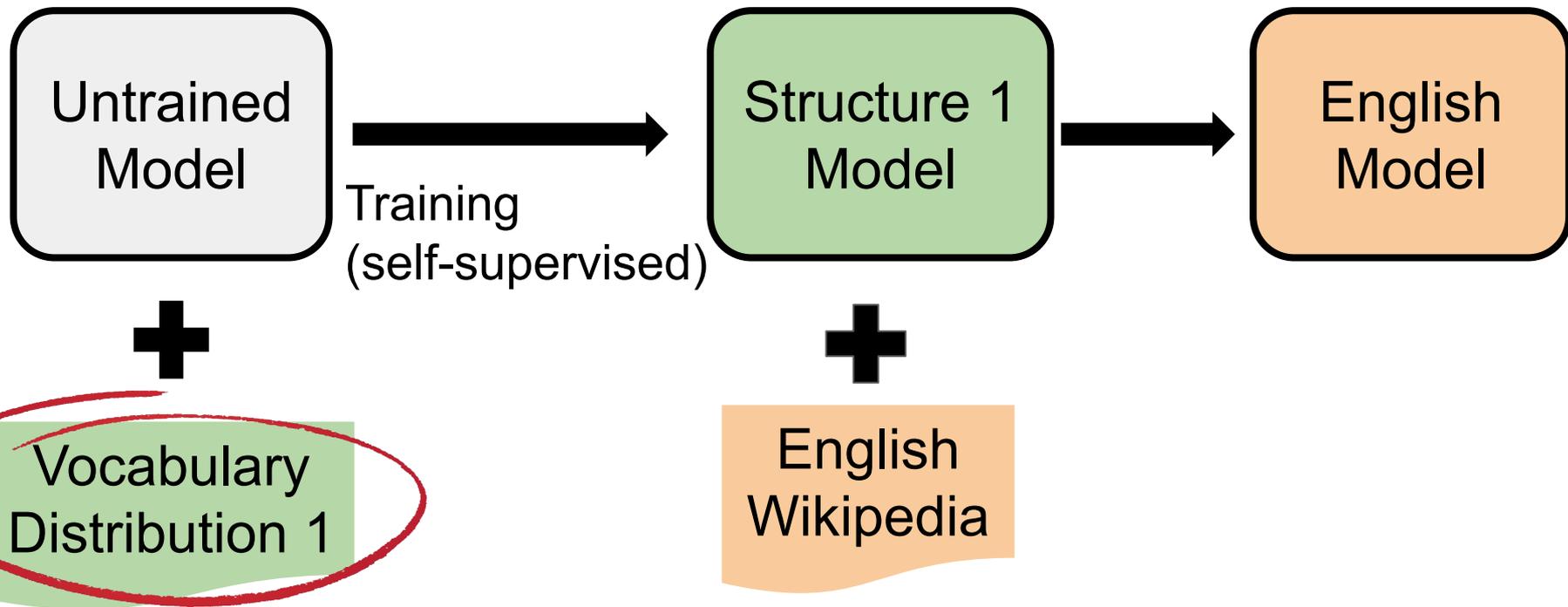
cat, dog ... book ... ameliorate ...

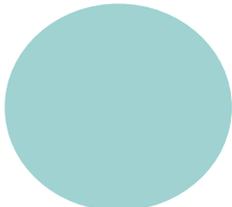


Vocabulary indices 0-50K



Structural Transfer: a testbed for linguistic structure hypotheses

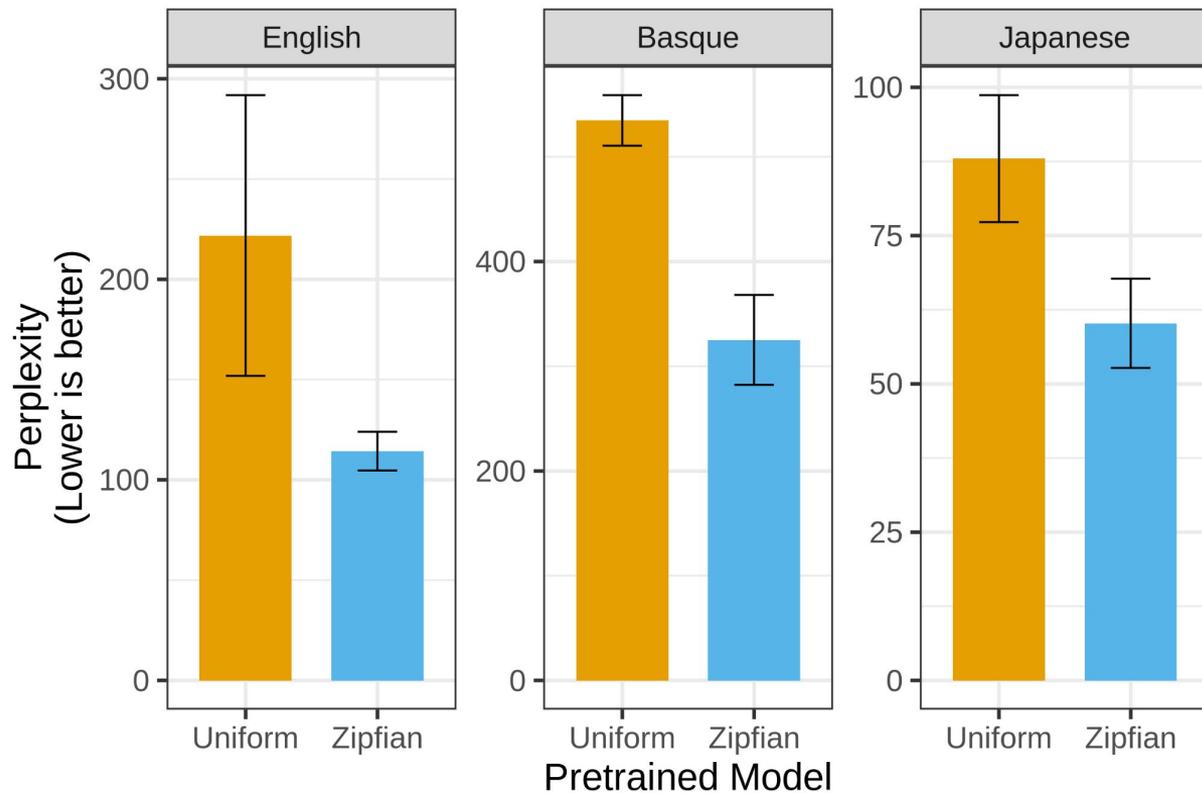




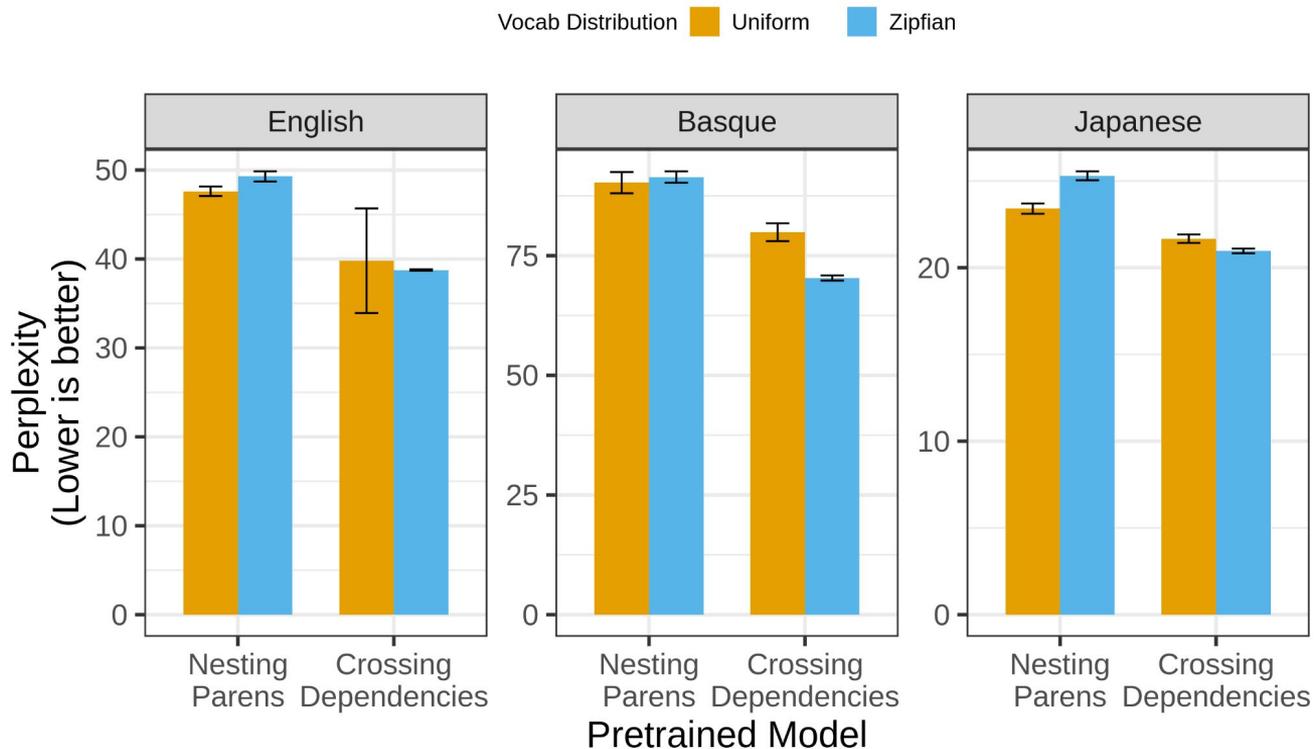
Does a Zipfian vocabulary distribution in pretraining have a **structural** effect?

- Even though we discard vocabulary information

Yes, Zipfian information is transferred



... but does not necessarily combine with structure



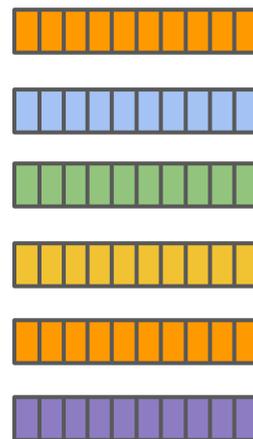
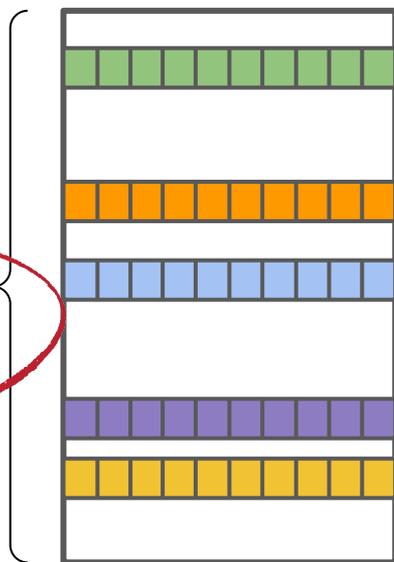
The role of vocabulary in transfer learning is an interesting problem

Vocabulary
embedding matrix

Lexical
representations

“The cat
sat on the
mat”

$|V| \approx$
10s of K



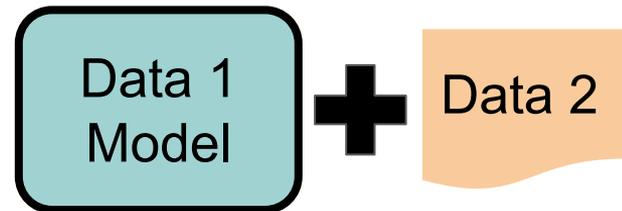
Really big
transformer
model

The role of vocabulary in transfer learning is an interesting problem

- A practical problem: without enough data, it's hard to see a word often enough to learn a good vector
- A puzzle: how is structural information separated between vocabulary matrix and model weights?
 - Vocabulary information like distribution can have structural effects

Transfer learning, language, and structure

- Transfer learning is a test bed for understanding structure in language learning



- Computational models of cognitive processes can't prove anything – but they serve as interesting hypotheses generators



- It's an exciting time: machine learning opens up new avenues for exploring questions in language



Thanks!