



LOL Lab

Language Of  
Learning Lab

Est. 2019

# Meaning Between the Lines: Using Multidimensional Linguistics and Machine- Learning to Predict Reading Comprehension

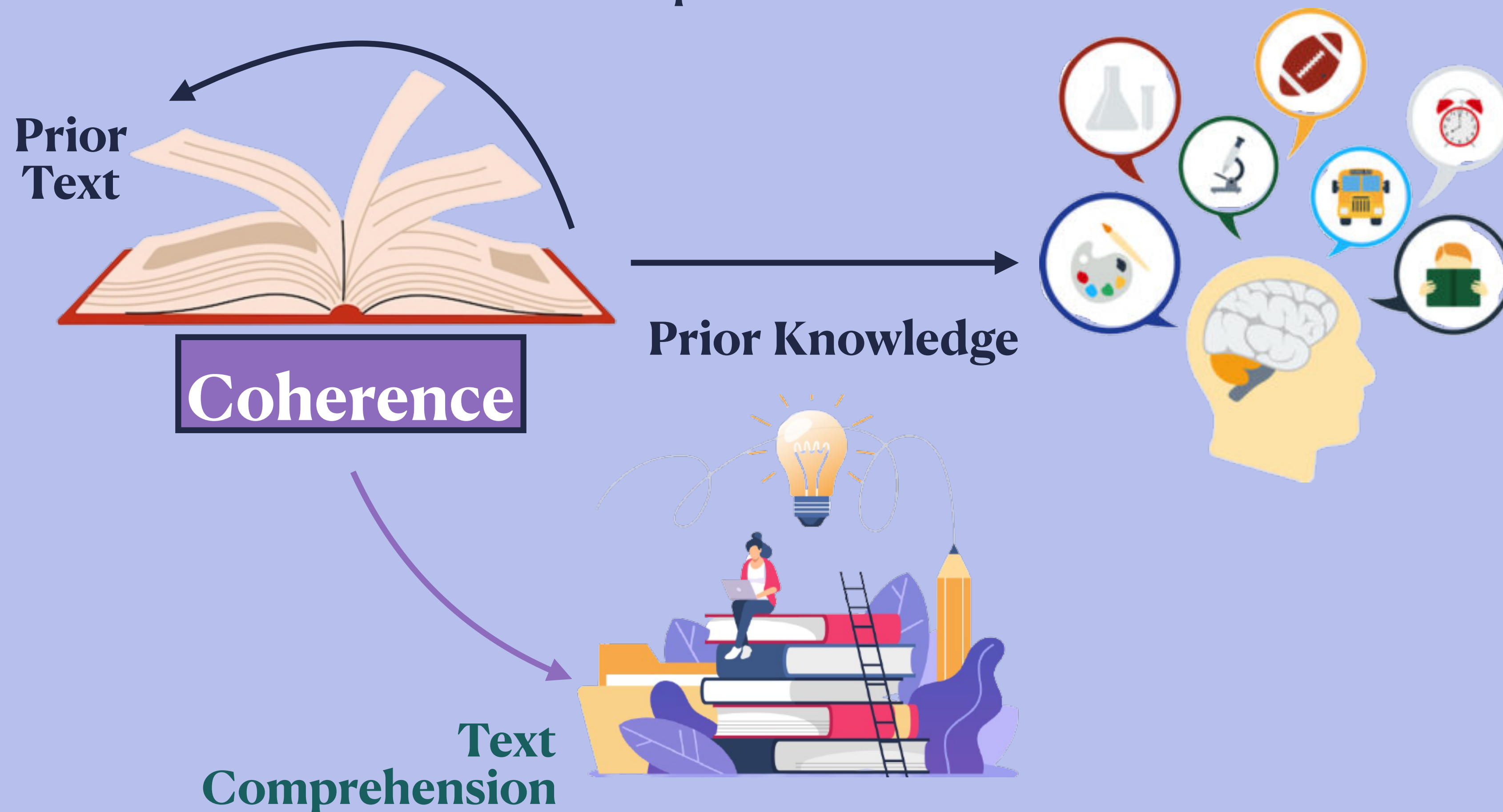
Lauren E. Flynn



November 2022

# successful text comprehension

when a reader has constructed a *coherent* & meaningful mental representation of a text



# coherence vs. cohesion

---

relational property between  
the reader and the text

how connected ideas are of  
readers' interpretations of text

necessary for comprehension

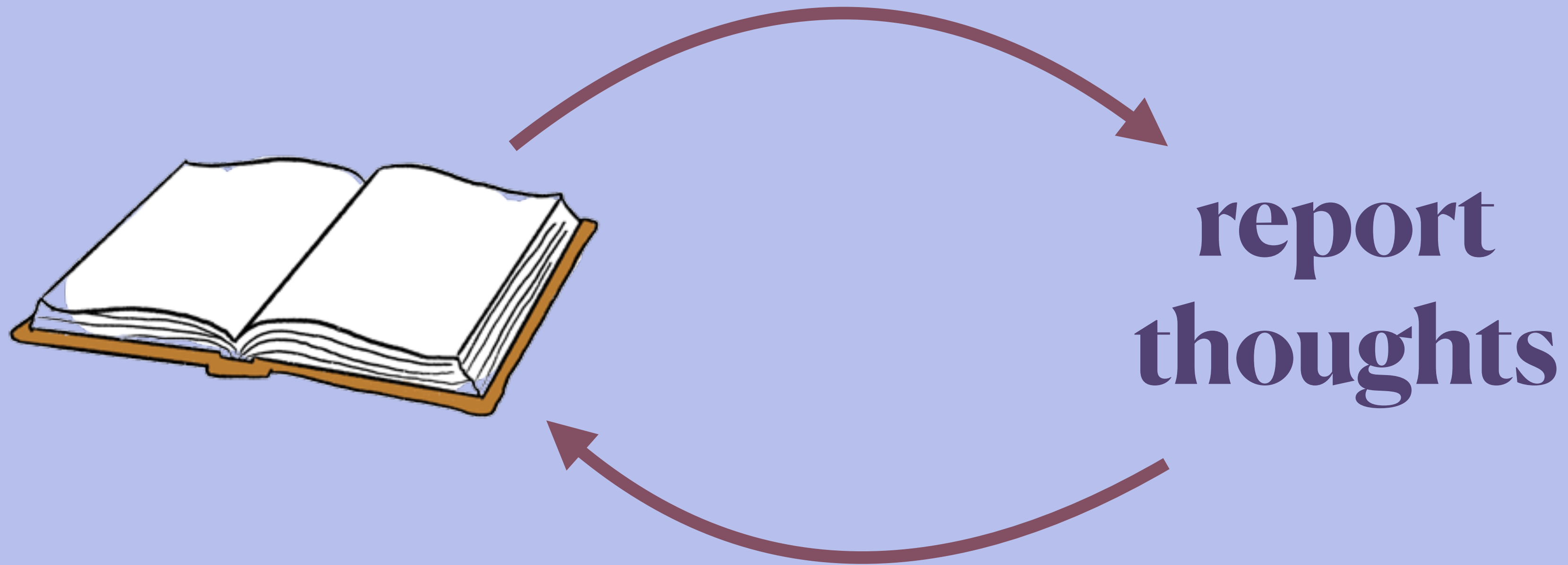
property of text

how connected words and  
sentences are in discourse

can aid in comprehension

# think-aloud procedures

examines coherence-building  
processes while reading



# think-aloud cohesion

depends on reader successfully developing  
coherent representations while reading

cohesion of think-alouds used to measure readers'  
coherence of text

positively associated with:

**reading  
skills**

**vocabulary  
knowledge**



# self-explanations

type of comprehension strategy used during think-alouds

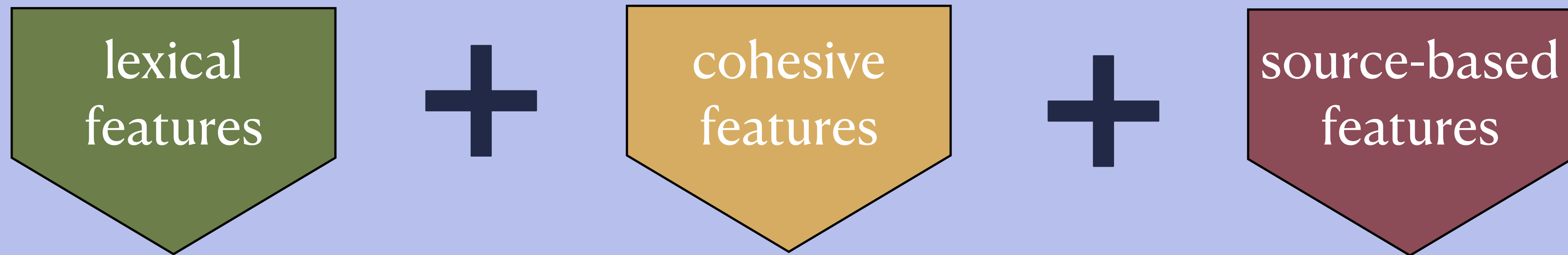
**explaining meaning of  
current text to one's self**

promotes inference-generation, deep comprehension, &  
problem solving

movement towards leveraging NLP to  
assess constructed responses

# current study

analyzes multidimensional linguistic features of self-explanations and machine-learning to predict readers' comprehension



using a diverse set of science texts



**methods**

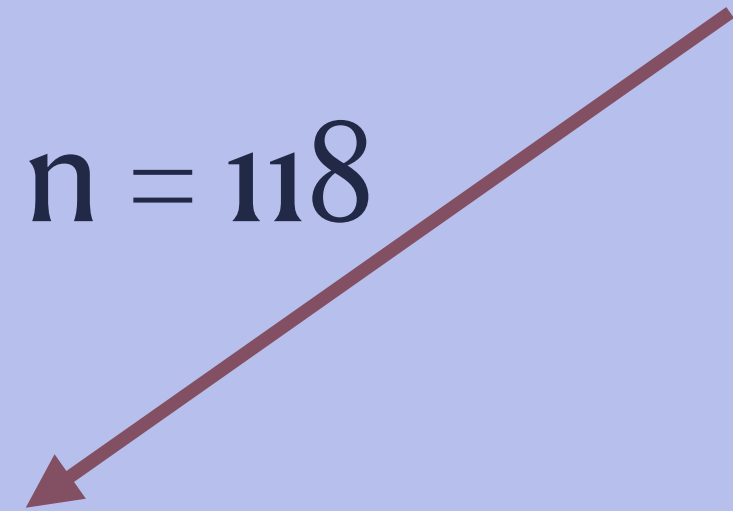


# methods



$N = 511$

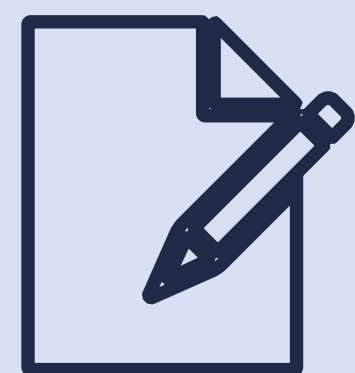
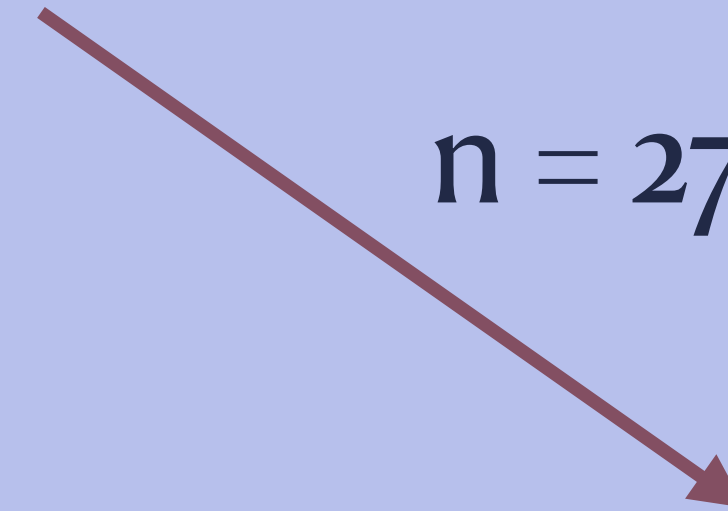
$n = 118$



$n = 115$



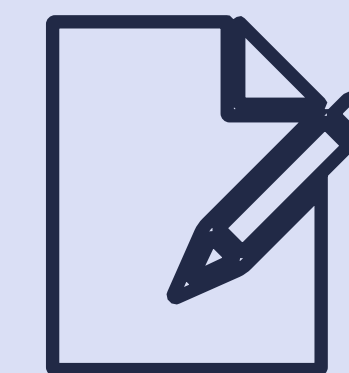
$n = 278$



shallow Q's  
deep Q's



shallow Q's  
deep Q's



shallow Q's  
deep Q's

# comprehension scoring

human raters scored the 8 open-ended comprehension questions for partial credit, ranging from 0-1 for both the textbase and bridging inferences scores

0.00	0.25	0.50	0.75	1.00
------	------	------	------	------

**textbase (shallow)**

0.75

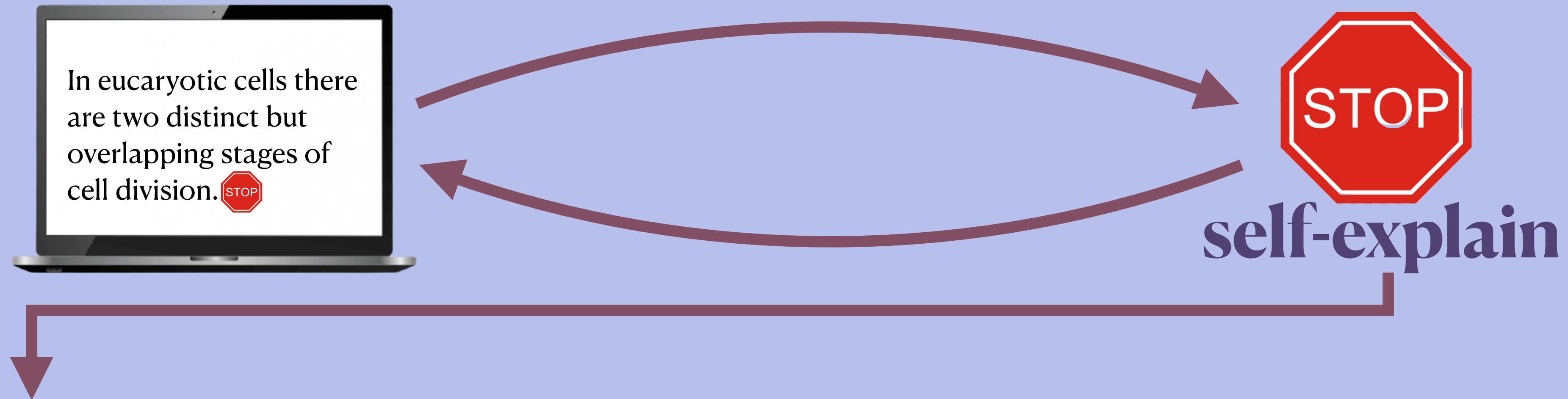
**bridging inferences (deep)**

0.25

$$\frac{0.75 + 0.25}{2} = 0.50$$

(total comprehension score)

# procedure



## lexical features

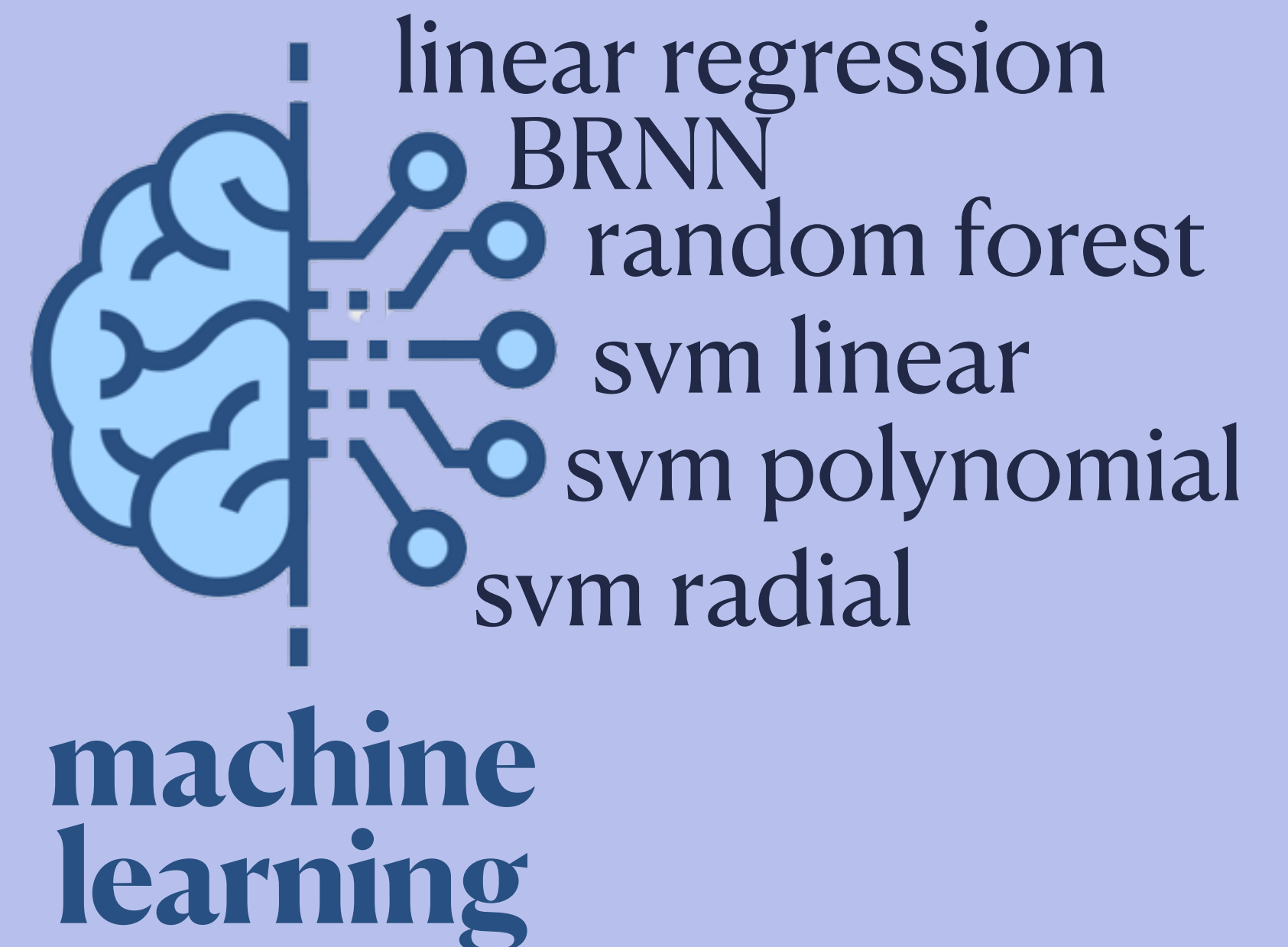
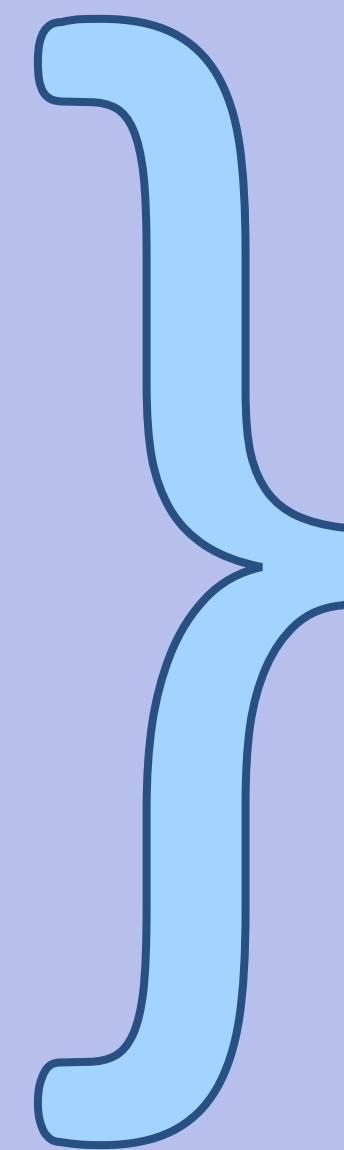
age of acquisition  
academic content word frequency  
concreteness of content words

## cohesive features

noun type-token-ratio  
Word2vec semantic similarity  
sentence linking

## source-based features

% noun type overlap  
% of lemma overlap  
total verb lemma overlap





# Natural Language Processing (NLP)

lexical features

cohesive features

source-based features

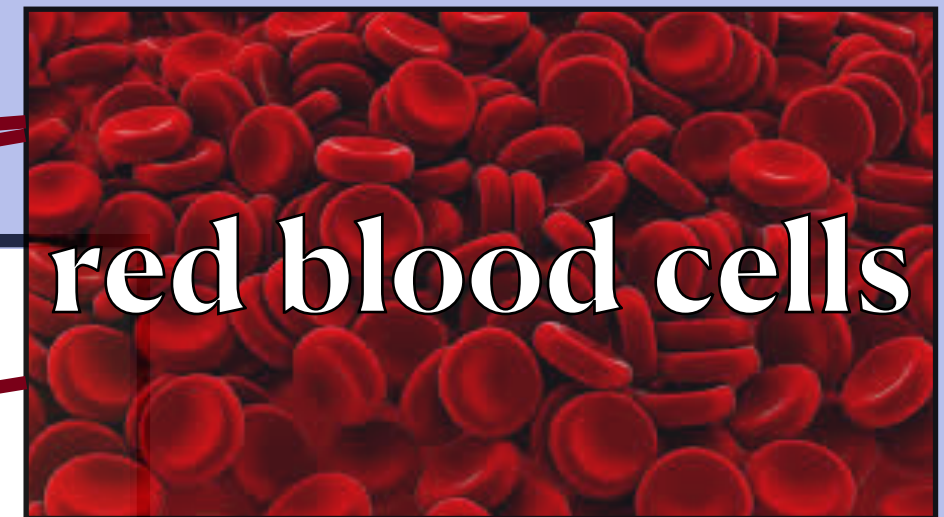
Red blood cells are a necessity for the body. They bring carbon dioxide to the cells of the body. The body then turns the oxygen into carbon dioxide. The red blood cells take the carbon dioxide and have it removed

.....

healthy red blood cells carry oxygen around your body. Sometimes not enough oxygen is transported. Not enough oxygen results in anemia

.....

Anemia is a condition where not enough oxygen gets into the body. Anemia can make a person feel tired and weak. One time, my doctor told me I was anemic and it made me feel really tired. I guess, anemia must have something to do with problems with your red blood cells. Is blood a cell?

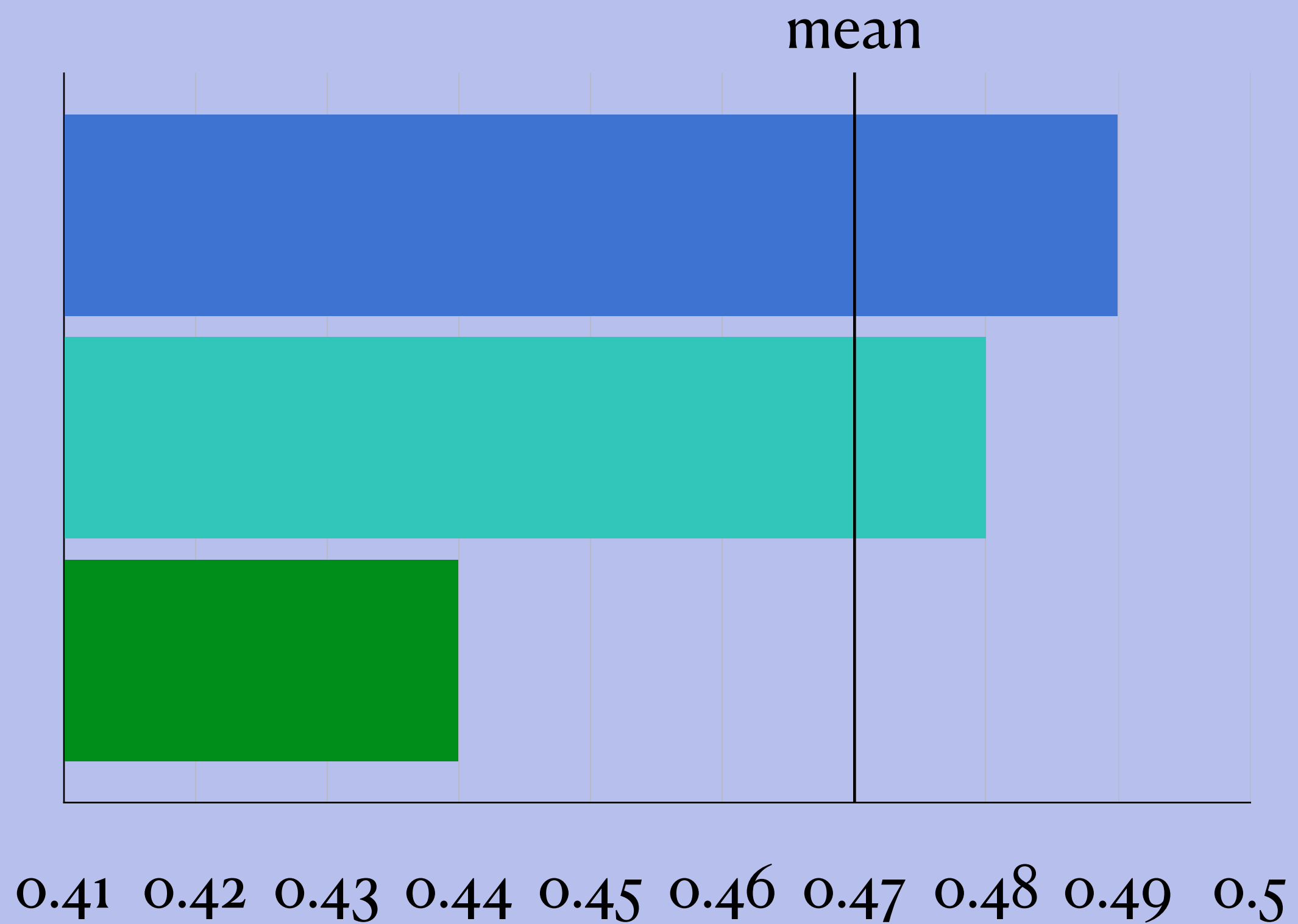


**results**

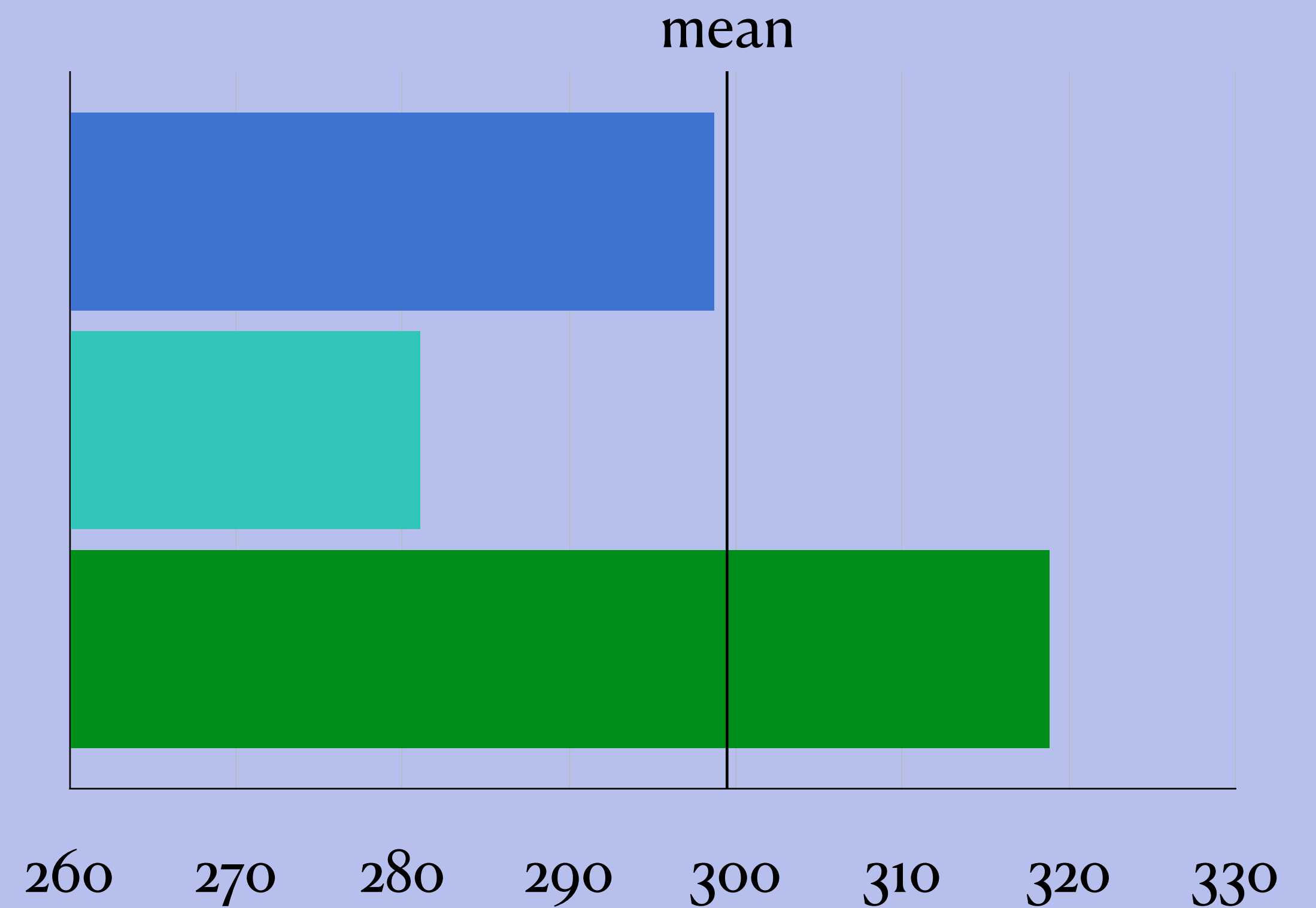


# results: descriptives

comprehension scores



# of words



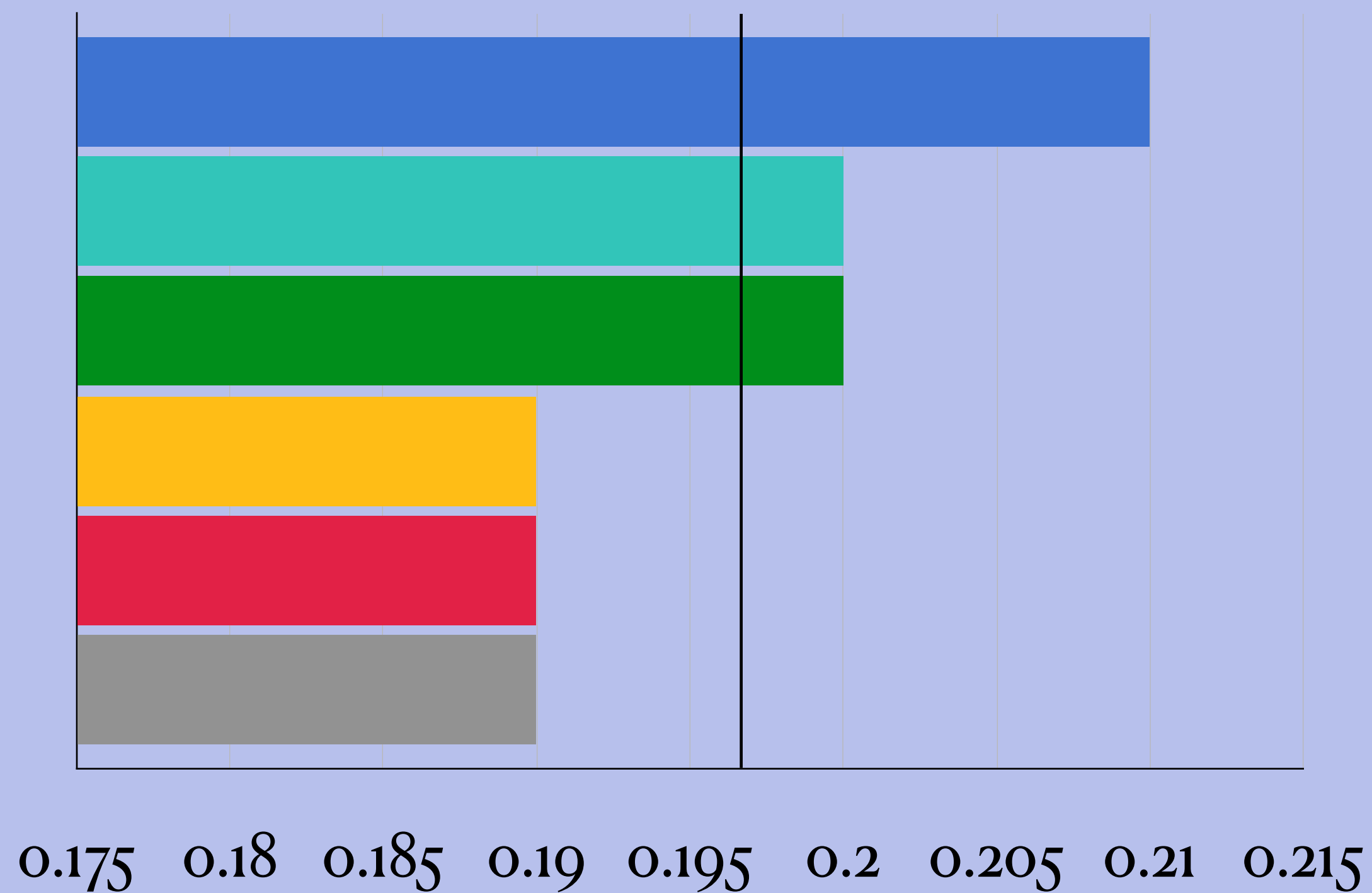
heart disease  
cell division

red blood cells

# results: algorithm performance

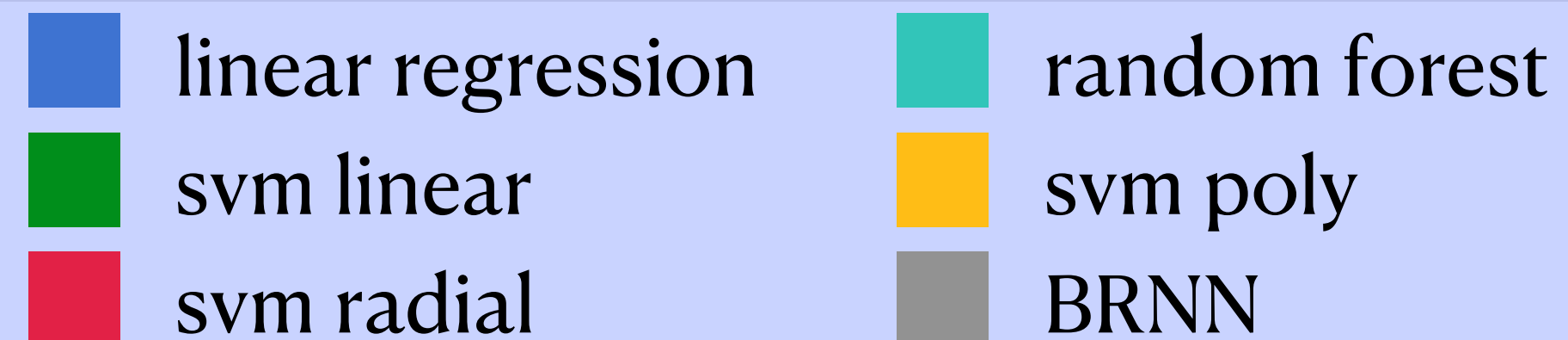
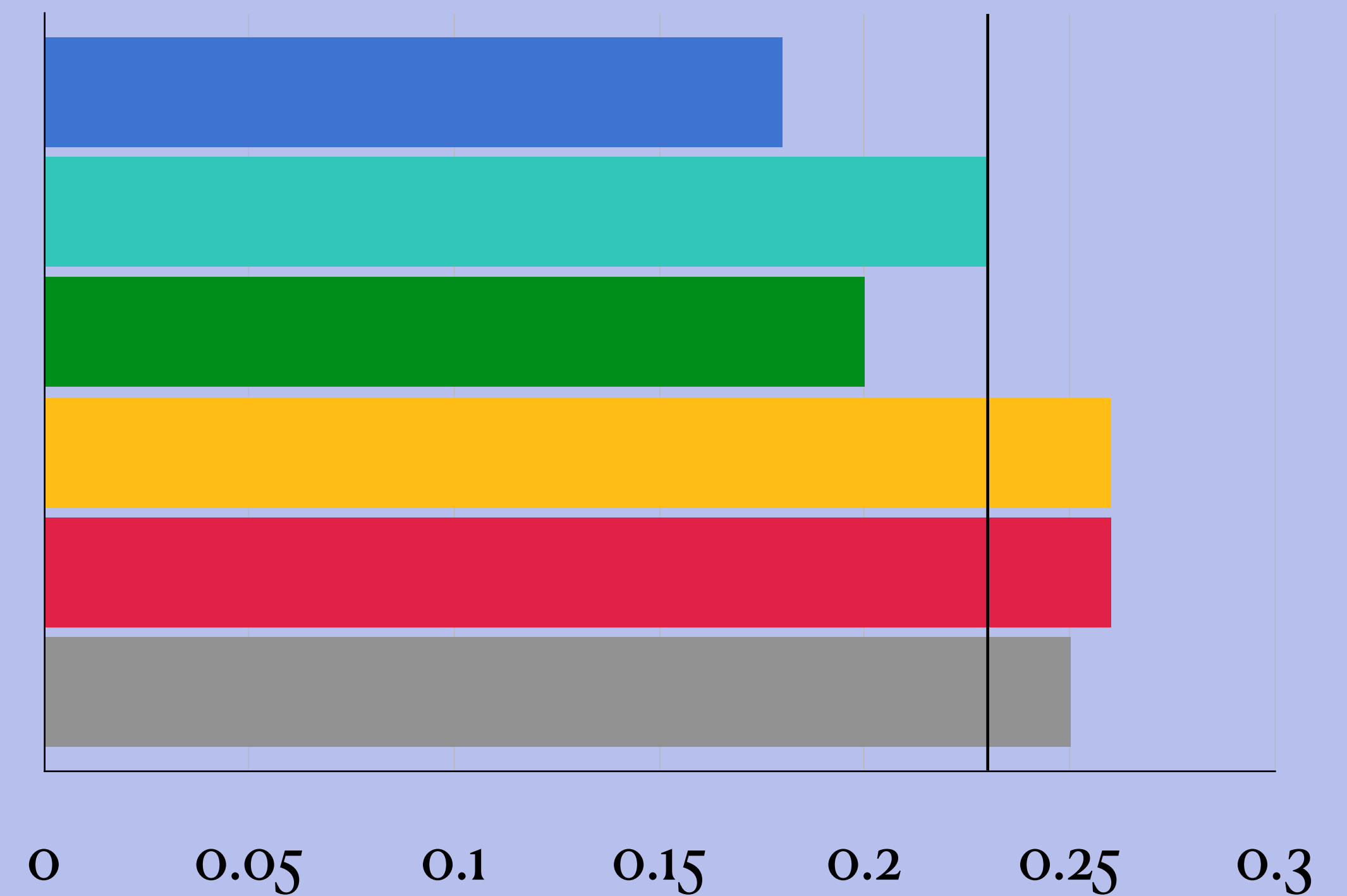
## RMSE

mean



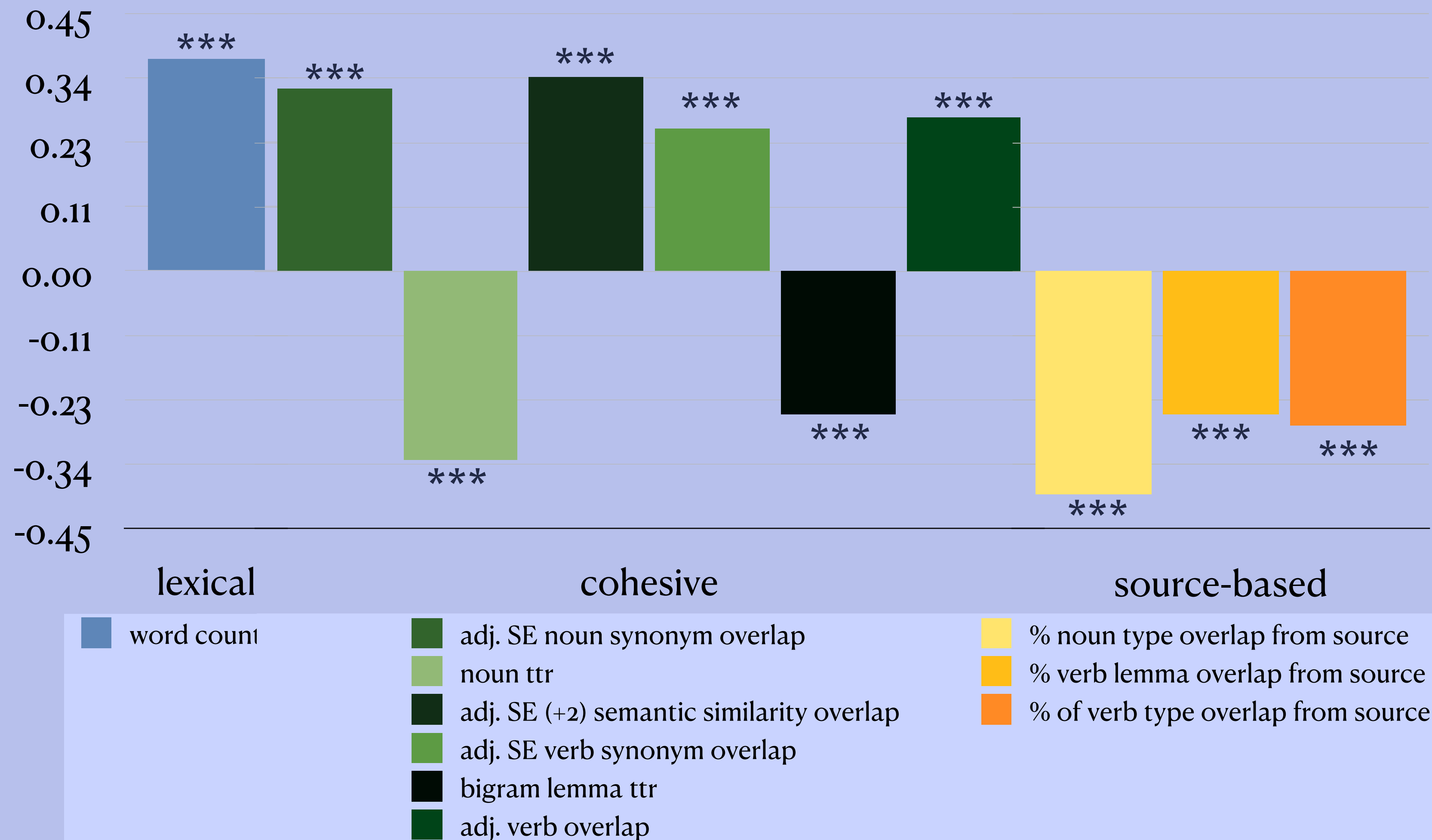
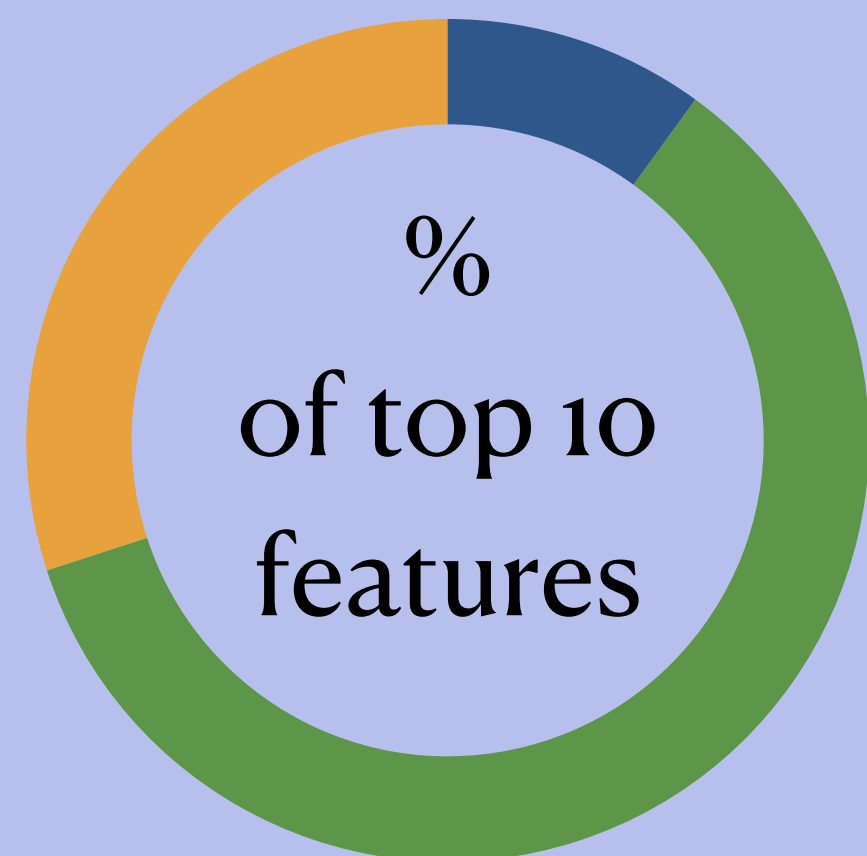
## R-squared

mean



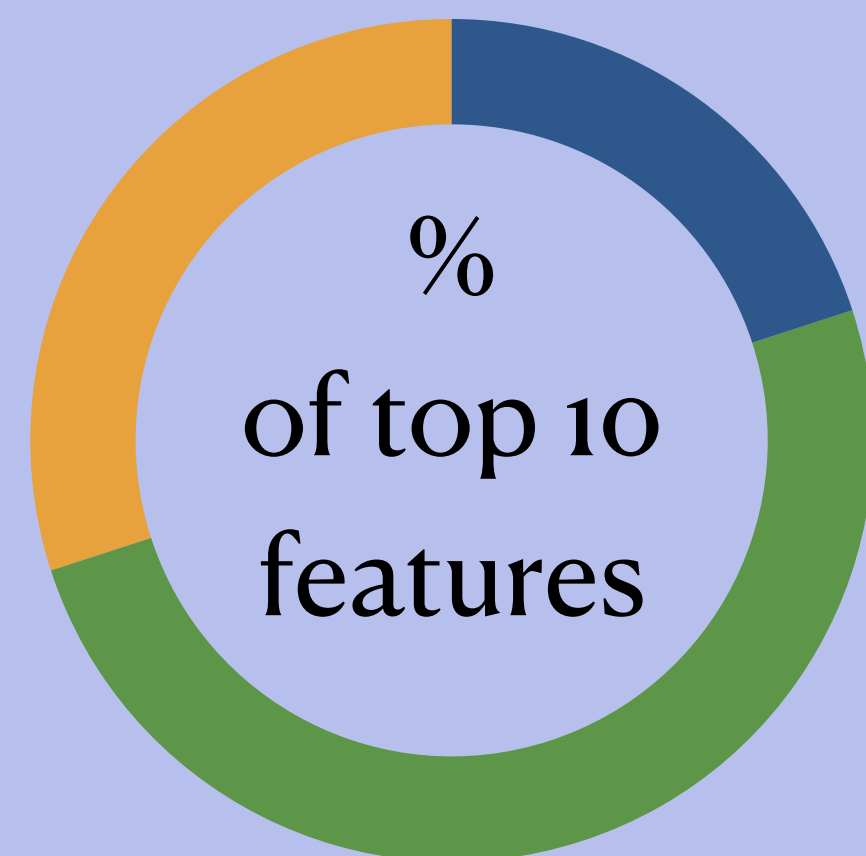
# results: top 10 correlated features

- lexical
- cohesive
- source-based



# results: model exploratory feature analysis

- lexical
- cohesive
- source-based

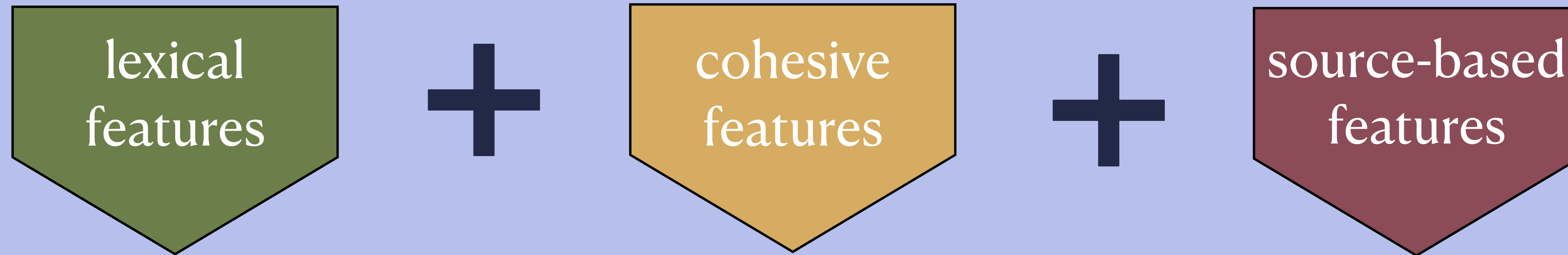


**discussion**



# summary

leveraged multidimensional linguistic features



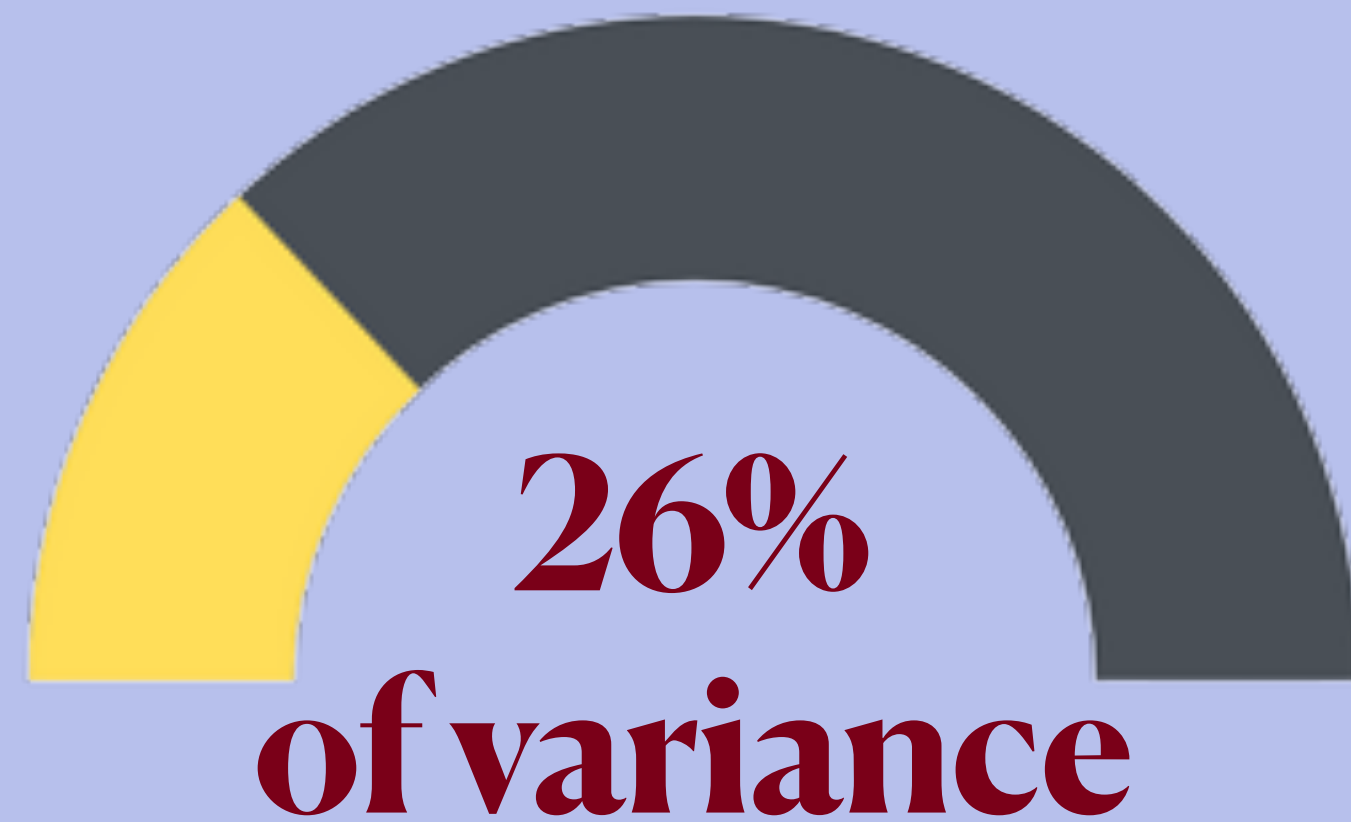
and machine-learning algorithms

	linear regression		random forest
	svm linear		svm poly
	svm radial		BRNN

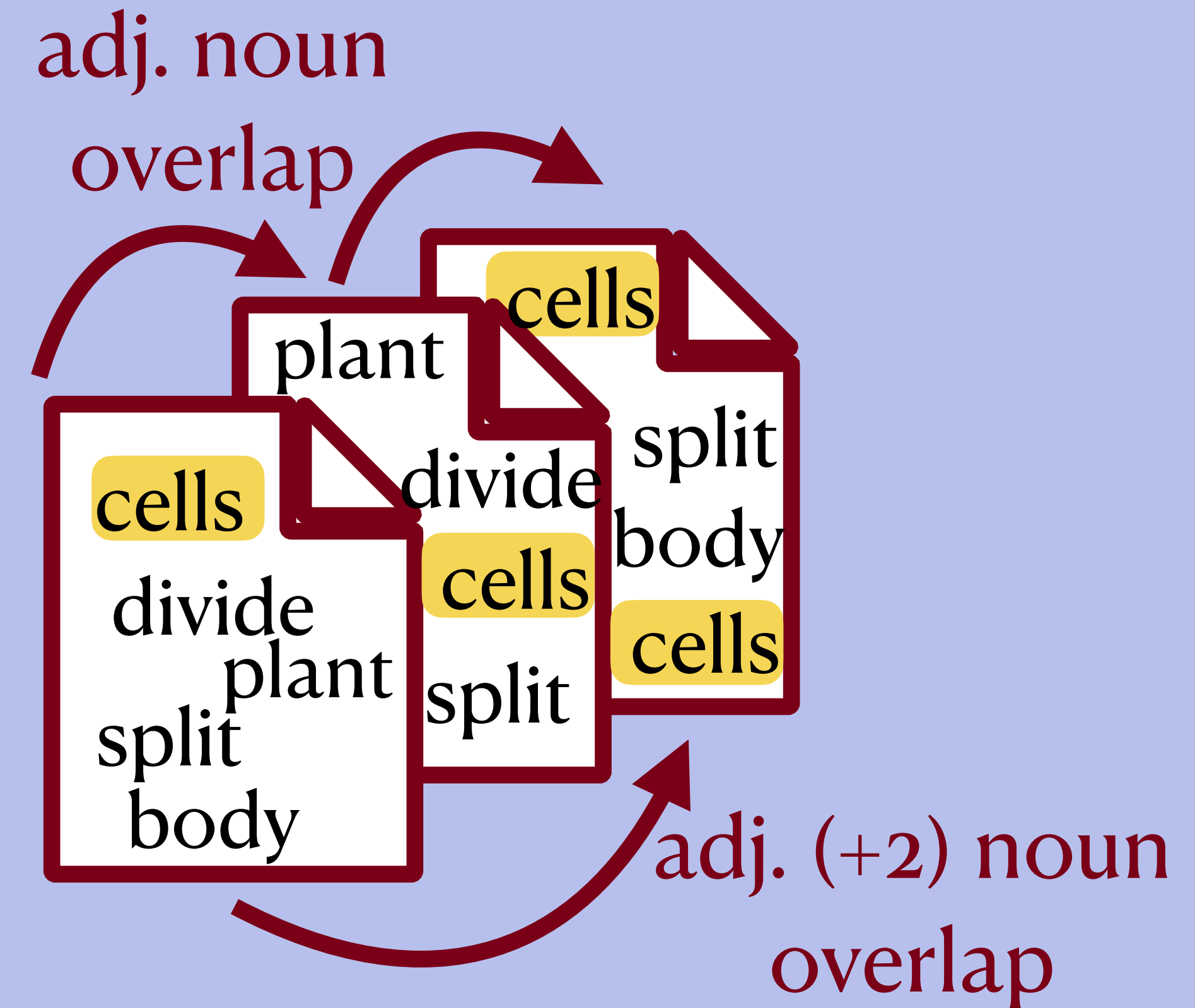
to predict overall comprehension of science texts



# study significance



svm radial algorithm was  
the best performing



noun overlap (cohesion & source-based) were most important to model performance



# moving forward

more fine-grained linguistic analyses to determine linguistic category most indicative of student comprehension

lexical  
features

**VS.**

cohesive  
features

**VS.**

source-based  
features

expand scope to other genres of text

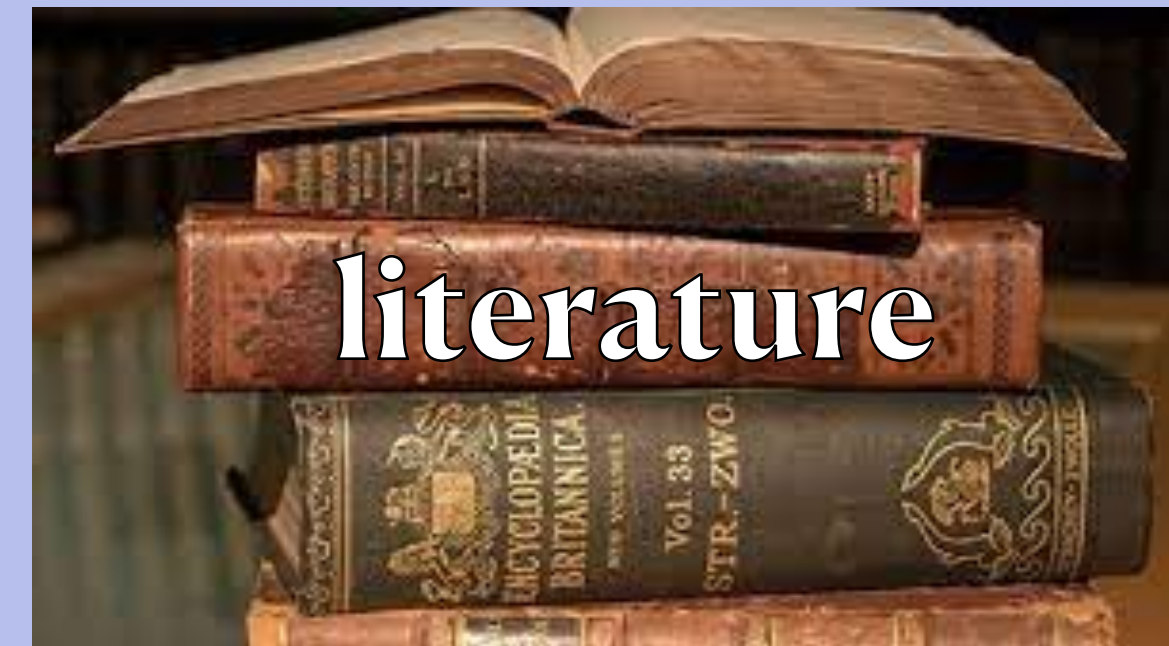
history



news



literature



# moving forward

more fine-grained analyses to assess linguistic feature importance to indications of shallow or deep comprehension



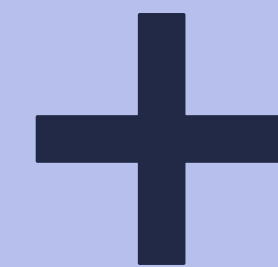
shallow processing (textbase-level comprehension)

deep processing (situation-model comprehension)

e.g., bridging inferences

examine additional linguistic dimensions

sentiment  
features



syntactic  
features

**questions?**

**email: [flynn598@umn.edu](mailto:flynn598@umn.edu)**