

Summer School

Large Language Models for Digital Humanities Research

Track 3: Using LLMs for Psycholinguistic Research

Part 2: Text Generation and Text Analysis

Job Schepens

job.schepens@uni.koeln.de

Contents

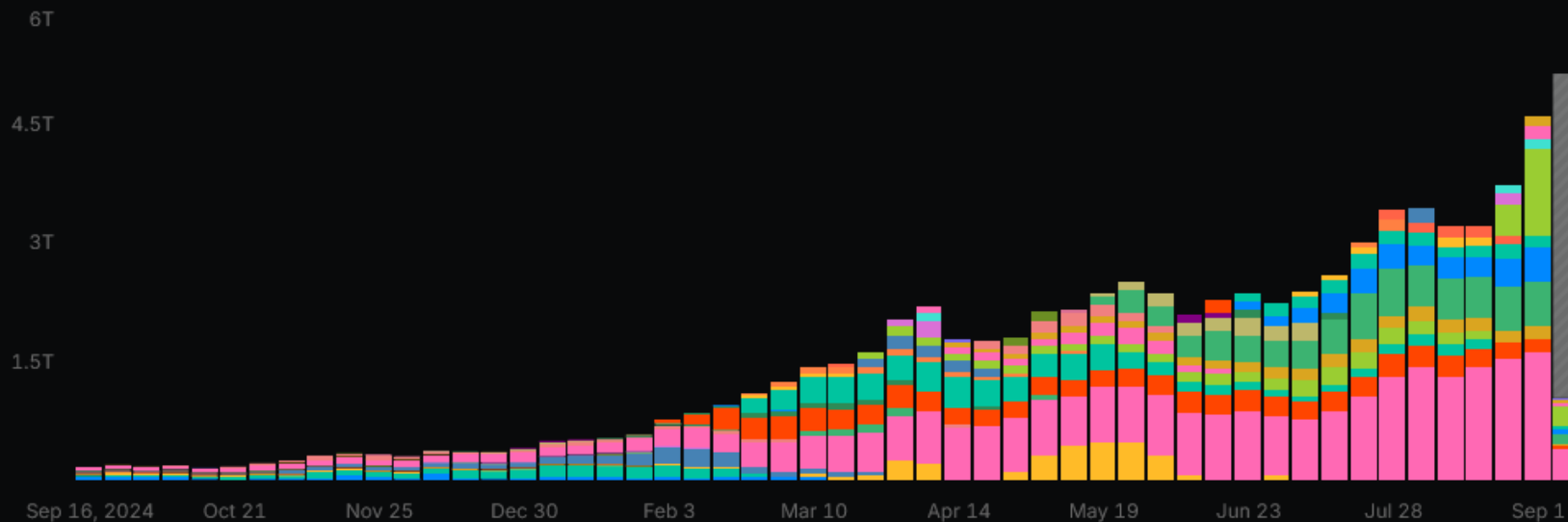
1. Current LLMs
2. Using LLMs for experimental stimulus generation
 - Two toy examples using vibe coding
3. Using LLMs for large-scale text generation
 - Schepens, Wołoszyn, Marx, & Gagl (accepted July 2025, MIT Open Mind).
4. Practical session using python notebooks:
 - **1st part:**
 - **Notebook 1: Experimenting with LLM-based corpus generation**
 - The repo also contains a few pre-generated 2m corpora (see /scripts folder in the repo)
 - Optional notebooks: Checking the corpus, extracting and formatting with metadata, merging data with behavioral data, comparing different frequency measures, comparing transformations
 - **2nd part:**
 - **Notebook 2: Validating predictors against human reading times**

<https://github.com/jobscapens/mlschool-text>

Leaderboard

Top this week ↕

Token usage across models on OpenRouter ⓘ



1.	Grok Code Fast 1 by x-ai	1.15T tokens ↑ 115%	6.	DeepSeek V3 0324 by deepseek	170B tokens ↑ 12%
2.	Claude Sonnet 4 by anthropic	564B tokens ↑ 3%	7.	DeepSeek V3.1 (free) by deepseek	149B tokens ↑ 540%
3.	Gemini 2.5 Flash by google	325B tokens ↓ 31%	8.	Gemini 2.5 Pro by google	148B tokens ↓ 18%
4.	Gemini 2.0 Flash by google	181B tokens ↓ 11%	9.	Qwen3 30B A3B by qwen	113B tokens ↓ 4%
5.	GPT-4.1 Mini by openai	173B tokens ↑ 387%	10.	GPT-5 by openai	85.4B tokens ↑ 15%

Llama 3.3 70B

"This model delivers similar performance to Llama 3.1 405B with cost effective inference that's feasible to run locally on common developer workstations."

405B

70B

December 2024

DeepSeek v3 for Christmas

685B, estimated training cost \$5.5m

February 2025

Mistral Small 3 (24B)

"Mistral Small 3 is on par with Llama 3.3 70B instruct, while being more than 3x faster on the same hardware."

~20GB, January 2025

March 2025

GPT 4.1 (1m tokens!)

gpt-4.1-nano

gpt-4.1-mini

gpt-4.1

o3 and o4-mini

o3

o4-mini

May

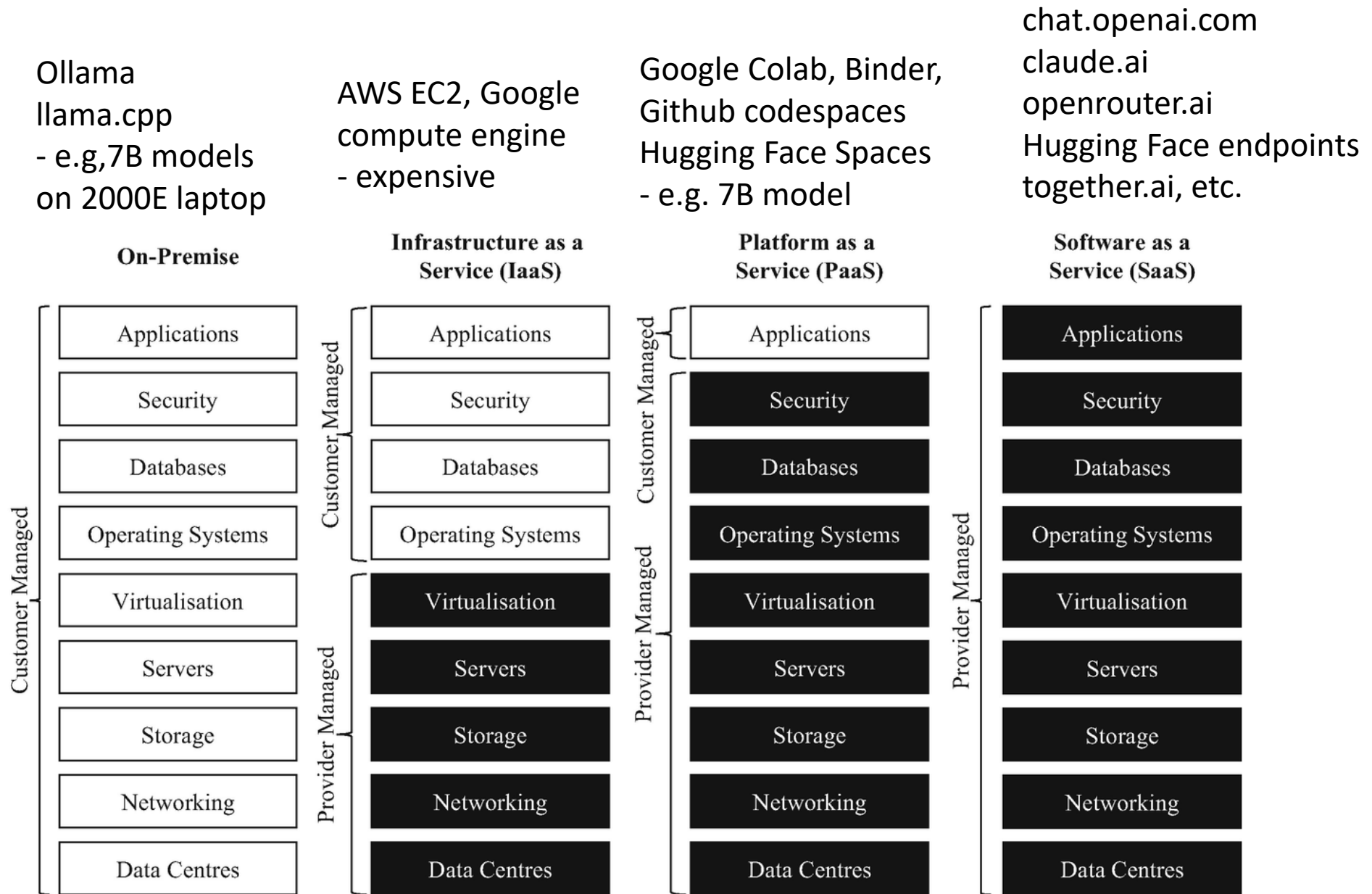
Claude Sonnet 4

Claude Opus 4

gemini-2.5-pro-
preview-05-06

Current developments in NLP using LLMs

- Many current developments...
 - Reasoning, agentic capabilities, larger context windows, multimodal integration, open weight models
- Many (immature) possibilities
 - As (safe?) **tools**: e.g. stimulus generation, automatic annotation, word frequency estimation
 - As (robust?) **models** (of what?): e.g. next word prediction / surprisal process and reading times / cloze probability
- But also many responsibilities (e.g. keynote by Elen le Foll)
 - Ethics, code of conduct, standards, best practices (transparency, human verification, etc.)
- LLMs across Language Research in Cologne
 - SFB1252 Research Data & Methods, SFB1252 Brown Bag Lunches, Reproducibilitea, DH Colloquium, UzK Data Steward Network, Informal discussion groups, etc.
 - Ongoing research: upcoming “LLMs for linguistic analysis” workshop in Cologne 24-25 November (more info soon)



https://en.wikipedia.org/wiki/Software_as_a_service#/media/File:Comparison_of_on-premise,_IaaS,_PaaS,_and_SaaS.png

LLMs for Experimental Stimulus Generation: Example 1

Generating sentences, controlling for syntax, frequency, and meaning:

- Syntactic complexity: Prompt engineering to generate sentences with specific syntactic structures (e.g., "Generate simple **SVO sentences**" vs. "Generate sentences with **embedded relative clauses**")
- Word frequency: Fine-tuning on frequency-controlled corpora to **maintain target word frequency ranges**
- Semantic content manipulation: Using **LLMs or embeddings**
- Multi-constraint generation: Simultaneous control of **multiple variables** (e.g., "Generate high-frequency words in complex syntactic structures about cooking")

```

1 # Precise, topic-aware prompt engineering
2 prompts = {
3     ('SVO', 'high', 'cooking'):
4         "Write a simple sentence about cooking using common words:",
5     ('embedded', 'low', 'AI research'):
6         "Generate a sentence about artificial intelligence research with academic",
7     ('direct_object', 'medium', 'social'):
8         "Create a sentence about social interaction with moderate vocabulary that"
9 }
10
11 # spaCy-based syntactic analysis for validation
12 def analyze_syntax(self, sentence: str) -> Dict:
13     doc = self.nlp(sentence)
14
15     # Detect embedded relative clauses
16     relative_clauses = [token for token in doc if token.dep_ == 'relcl']
17
18     # Count direct objects
19     direct_objects = [token for token in doc if token.dep_ == 'dobj']

```

```

1 def generate_multiple_candidates(self, prompt: str, num_candidates: int):
2     candidates = []
3     inputs = self.tokenizer.encode(prompt, return_tensors='pt')
4
5     for _ in range(num_candidates):
6         outputs = self.model.generate(
7             inputs,
8             max_length=len(inputs[0]) + 25,
9             temperature=0.9,          # Controlled randomness
10            top_p=0.9,                 # Nucleus sampling
11            repetition_penalty=1.1,    # Avoid repetition
12            do_sample=True
13        )
14
15        # Extract and clean first sentence
16        generated = self.tokenizer.decode(outputs[0], skip_special_tokens=True)
17        sentence = generated[len(prompt):].split('.')[0] + '.'
18        candidates.append(sentence)

```

Generated Sentence	Syntactic Target	Frequency Target	Topic	Syntax Match	Frequency Match	Topic Score	All Constraints Met	Attempts Needed	Dependency Depth	Length
SVO + HIGH FREQUENCY + COOKING										
The time is coming.	SVO	HIGH	cooking and foo...	✓ YES	✓ YES	0.722	🏆 PERFECT	96	2	5
A person's name is your number.	SVO	HIGH	cooking and foo...	✓ YES	✓ YES	0.714	🏆 PERFECT	139	3	8
The chicken is cooked very well.	SVO	HIGH	cooking and foo...	✓ YES	✗ NO	0.795	⚠️ PARTIAL	117	2	7
The recipe includes vegetables.	SVO	HIGH	cooking and foo...	✓ YES	✗ NO	0.792	⚠️ PARTIAL	78	2	5
a chicken curry or something.	SVO	HIGH	cooking and foo...	✓ YES	✗ NO	0.786	⚠️ PARTIAL	68	1	7
EMBEDDED RELATIVE + LOW FREQUENCY + AI RESEARCH										
Machine Learning (MLA) is the field where computer science becomes more interesting as it progresses through each area of study.	EMBEDDED_RELATIVE	LOW	artificial inte...	✓ YES	✓ YES	0.875	🏆 PERFECT	11	7	23
say in the example where AI may be used as well.	EMBEDDED_RELATIVE	LOW	artificial inte...	✓ YES	✓ YES	0.795	🏆 PERFECT	24	5	12
It is conceivable to design something from the data it can learn.	EMBEDDED_RELATIVE	LOW	artificial inte...	✓ YES	✓ YES	0.810	🏆 PERFECT	29	6	13
The next thing I want to do is talk about the future.	EMBEDDED_RELATIVE	LOW	artificial inte...	✓ YES	✓ YES	0.780	🏆 PERFECT	32	4	13
In order to learn how the human mind works we need an environment that can understand our emotional state.	EMBEDDED_RELATIVE	LOW	artificial inte...	✓ YES	✓ YES	0.816	🏆 PERFECT	33	6	20
DIRECT OBJECTS + MEDIUM FREQUENCY + SOCIAL										
I think this is why we love each other, or the relationship I live in will change.	DIRECT_OBJECTS	MEDIUM	social interact...	✓ YES	✓ YES	0.747	🏆 PERFECT	3	5	21
The word emotional is often used to describe the process by which people feel themselves felt.	DIRECT_OBJECTS	MEDIUM	social interact...	✓ YES	✓ YES	0.774	🏆 PERFECT	5	5	19
The most successful people are those who have an interest in what they're doing.	DIRECT_OBJECTS	MEDIUM	social interact...	✓ YES	✓ YES	0.778	🏆 PERFECT	7	6	16
We want to write as many words per minute in an hour or two of time.	DIRECT_OBJECTS	MEDIUM	social interact...	✓ YES	✓ YES	0.723	🏆 PERFECT	9	6	17
I have an opinion on this person, He is very nice or simply You are extremely polite.	DIRECT_OBJECTS	MEDIUM	social interact...	✓ YES	✓ YES	0.695	🏆 PERFECT	11	4	24

LLMs for Experimental Stimulus Generation: Example 2

Generating minimal pairs

- Phonological: Generating words using an **encoder model (e.g. BERT)** to control semantic plausibility and filter for similarity in onset, length, stress, tone, vowel quality, voicing, aspiration, ...
- Morphological: Generation of **inflectional and derivational minimal pairs** (e.g., walk/walked, happy/happiness)
- Syntactic: Generating sentences that differ only in **target syntactic structures**
- Cross-modal minimal pairs: Generating **text-image pairs** with controlled linguistic-visual correspondences

Pipeline:

1. Template Creation

```
"The [MASK] was parked outside"
```

2. BERT Prediction

```
1 candidates = model.predict(template)
2 # → car, van, truck, bmw, cab...
```

3. Phonological Filtering

```
1 filter_by_contrast_type(
2     target="car",
3     candidates=candidates,
4     type="coda"
5 )
```

Controls:

- Edit distance = 1
- Position-specific contrasts
- Semantic plausibility
- Frequency matching

Context 1: “The red [MASK] was parked outside”

- car → cab (prob: 0.018)
- car → cat (prob: 0.003)

Context 2:
“She drove her [MASK] to work today”

- car → cat (prob: 0.001)

Context 3: “The [MASK] needs new tires”

- car → cat (prob: 0.002)
- car → cart (prob: 0.001)
- car → cab (prob: 0.001)

Generated: 6 minimal pairs across 3 contexts

Interim Summary: LLMs as tools (e.g. for stimulus generation)

Strengths

Weaknesses

Interim Summary: LLMs as tools (e.g. for stimulus generation)

Strengths

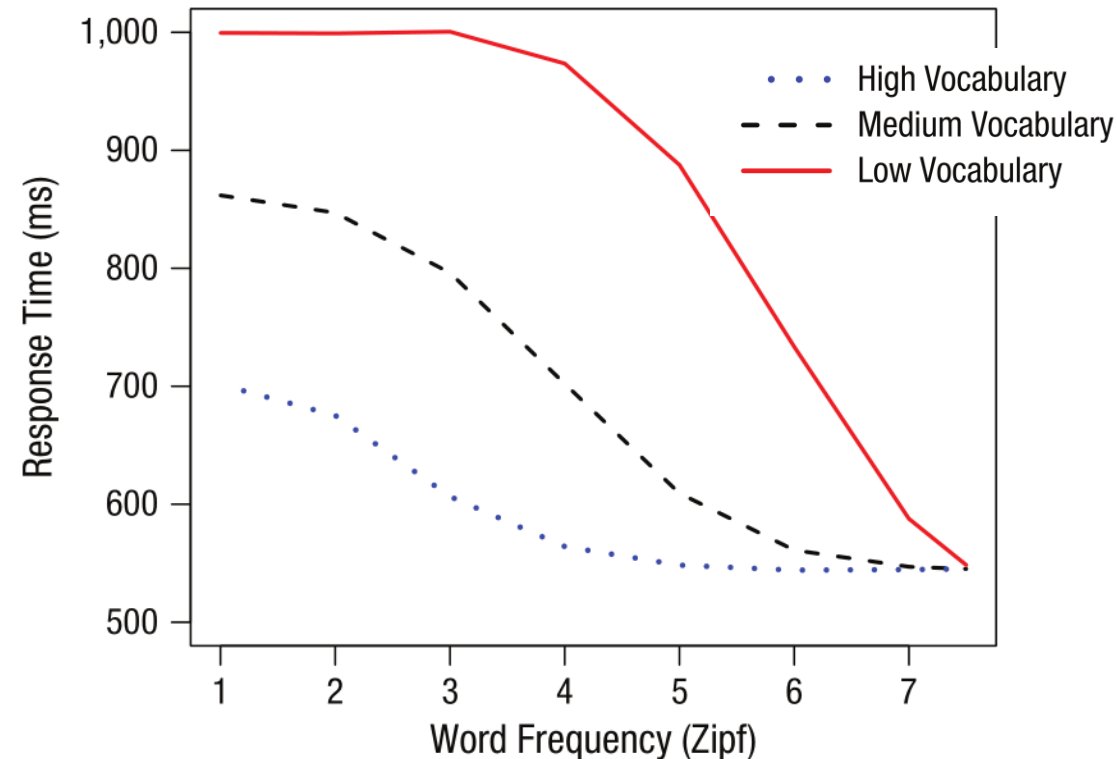
- **Scalability:** Generate many stimuli
- **Consistency and transparency:** Formalized criteria
- **Flexibility:** Easy to regenerate with different criteria
- **Reproducibility:** Documented and version controlled implementation

Weaknesses

- **Quality :** Biases due to LLM training data, fine-tuning, architecture, etc.
- **Lack of Domain Expertise:** Not trained on reasoning about specific linguistic issues
- **Reproducibility:** Due to LLM's stochastic nature, re-generating the code likely results in different stimuli
- **Ethical Concerns:** Possibly generating harmful content
- **Verification:** Human-in-the-loop
- **Black Box:** No mechanistic interpretation possible

Building linguistic corpora for word frequency estimation

- Word frequency is a strong behavioral correlate in visual word recognition paradigms
- Word frequency: counting word occurrences in a corpus
- Corpora used for estimating word frequency are based on text from:
 - Books, Newspapers DWDS, Heisters et al., 2011
 - Subtitles SUBTLEX, Brysbaert et al., 2011
 - German children's books childLex, Schroeder et al. 2015
 - Text generated by LLMs? Schepens et al., PsyArXiv



The Word Frequency Effect in Word Processing: An Updated Review

Marc Brysbaert¹, Pawel Mandera¹, and Emmanuel Keuleers²

¹Department of Experimental Psychology, Ghent University, and ²Department of Communication and Information Sciences, Tilburg University

Current Directions in Psychological Science
1–6
© The Author(s) 2017
Reprints and permissions:
sagepub.com/journalsPermissions.nav
DOI: 10.1177/0963721417727521
www.psychologicalscience.org/CDPS
SAGE

Two aims

Aim 1: **Compare word frequency** based on text generated by LLMs vs. text written by humans.

- Human text: ~10 million words in existing corpus of children's books (**ChildLEX**; Schroeder et al., 2015)
- Measures: correlation, number and percentage of shared words, lexical richness, Zipfs law, etc.

Aim 2: **Compare the estimated word frequency effect** on response times for LLM vs childLex word frequency

- Lexical decision response times for grade 1-6 and young and old adults (**DeveL**; Schröter & Schroeder, 2017)
- Measures: Improvement in model fit (AIC) of linear regression models, control for AoA, OLD20, word length

Generating 9 corpora (“conditions”)

GPT,
DeepSeek,
LLama

“Kinder” vs. “Erwachsene”

```
prompt=[  
  {  
    "role": "system",  
    "content": "4000 Wörter zu "  
    + titel  
    + " auf Deutsch geschrieben"  
    + " für Kinder im Alter "  
    + age_range  
  }  
]
```

Continue until
4000 words

```
openai.ChatCompletion.create(  
  model="gpt-3.5-turbo",  
  messages=prompt,  
  temperature=0.5,  
  max_tokens=4000,  
  n=4,  
  stop=None,  
  frequency_penalty=0,  
  presence_penalty=0
```

0.5 vs. 0.7

Generating 9 corpora

1. 1 corpus: GPT 3.5
2. 2x2 corpora: 2 temperatures (low, high) and 2 target audiences (child-directed, adult-directed prompt)
3. 2x2 corpora: 2 open weight models (DeepSeek V1, Llama 3.3 70B) and 2 text lengths (short, long)

GPT,
DeepSeek,
LLama

“Kinder” vs. “Erwachsene”

```
prompt=[  
  {  
    "role": "system",  
    "content": "4000 Wörter zu "  
    + titel  
    + " auf Deutsch geschrieben"  
    + " für Kinder im Alter "  
    + age_range  
  }  
]
```

Continue until
4000 words

```
openai.ChatCompletion.create(  
  model="gpt-3.5-turbo",  
  messages=prompt,  
  temperature=0.5,  
  max_tokens=4000,  
  n=4,  
  stop=None,  
  frequency_penalty=0,  
  presence_penalty=0
```

0.5 vs. 0.7

Example text (GPT 3.5)

Generated text: Das rote Phantom

Es war einmal ein kleiner **Junge namens Max**. Max war sieben Jahre alt und lebte in einer kleinen Stadt am Rande des Waldes. Er liebte es, draußen zu spielen und **Abenteuer** zu erleben.

Eines Tages, als Max im Wald spielte, hörte er ein seltsames Geräusch. Es war ein leises Knurren, das immer lauter wurde. Max drehte sich um und sah ein rotes Phantom auf ihn zukommen. Es hatte riesige Augen und scharfe Krallen.

Max war zuerst erschrocken, aber dann erkannte er, dass das rote Phantom nur ein kleiner Hund war. Der Hund war verängstigt und hungrig. Max beschloss, ihm zu helfen.

Klappentext

Kim, Franzi und Marie sind Die drei !!! . Mutig und clever ermitteln die drei Freundinnen und sind jedem Fall gewachsen.

Ein Phantom geht um! Nachts stellt ein **vermummter Motorradfahrer** auf seiner roten Maschine Verbrecher und vereitelt ihre Straftaten. Als ein dreister Dieb Franzis Rucksack klaut, ist das rote Phantom sofort zu Stelle. Der Ehrgeiz der drei Detektivinnen ist geweckt – wer ist das rote Phantom? Bei einer wilden Verfolgungsjagd durch die Nacht kommen Kim, Franzi und Marie hinter das Geheimnis und erleben eine erstaunliche Überraschung.



Example text (Llama 3.3 70B)

Generated text: Mein erster Schultag und der Eisbär-Schreck

Es war ein sonniger Morgen im September. Ich, Tim, war so aufgeregt, dass ich heute meinen ersten Schultag hatte. Meine Mutter half mir, mein neues Schulranzen zu packen, und mein Vater machte mir ein leckeres Frühstück.

Als wir zur Schule kamen, sah ich viele neue Gesichter. Ich fühlte mich ein bisschen schüchtern, aber meine Mutter gab mir einen großen Kuss und sagte: "Alles wird gut, Tim. Du wirst viele neue Freunde finden."

Ich ging in meine Klasse und sah, dass alle anderen Kinder auch ein bisschen nervös waren. Unser Lehrer, Herr Müller, stellte sich vor und erzählte uns, dass wir heute viele spannende Dinge lernen würden.

Klappentext

Ein fröhliches Mut-mach-Buch für den ersten Schultag

Ben freut sich riesig: Endlich kommt er in die Schule! Aber er hat auch ganz schön Lampenfieber, schließlich wird nach den Ferien alles anders sein: Er muss Plüschbärin Sardine zum ersten Mal alleine lassen! Und was wird aus seiner Schultüte – schließlich kann Mama doch gar nicht basteln! Wie gut, dass er Florence kennenlernt, die Bens beste Freundin wird. Denn wenn man zusammen in die Schule kommt, kann gar nicht mehr viel schiefgehen ...

- Mit farbenfrohen Illustrationen von Heike Wiechmann
- Einfühlsame Schilderung der aufregenden Zeit vor dem ersten Schultag
- Ein Vorlesebuch für alle, die auf den ersten Schultag warten



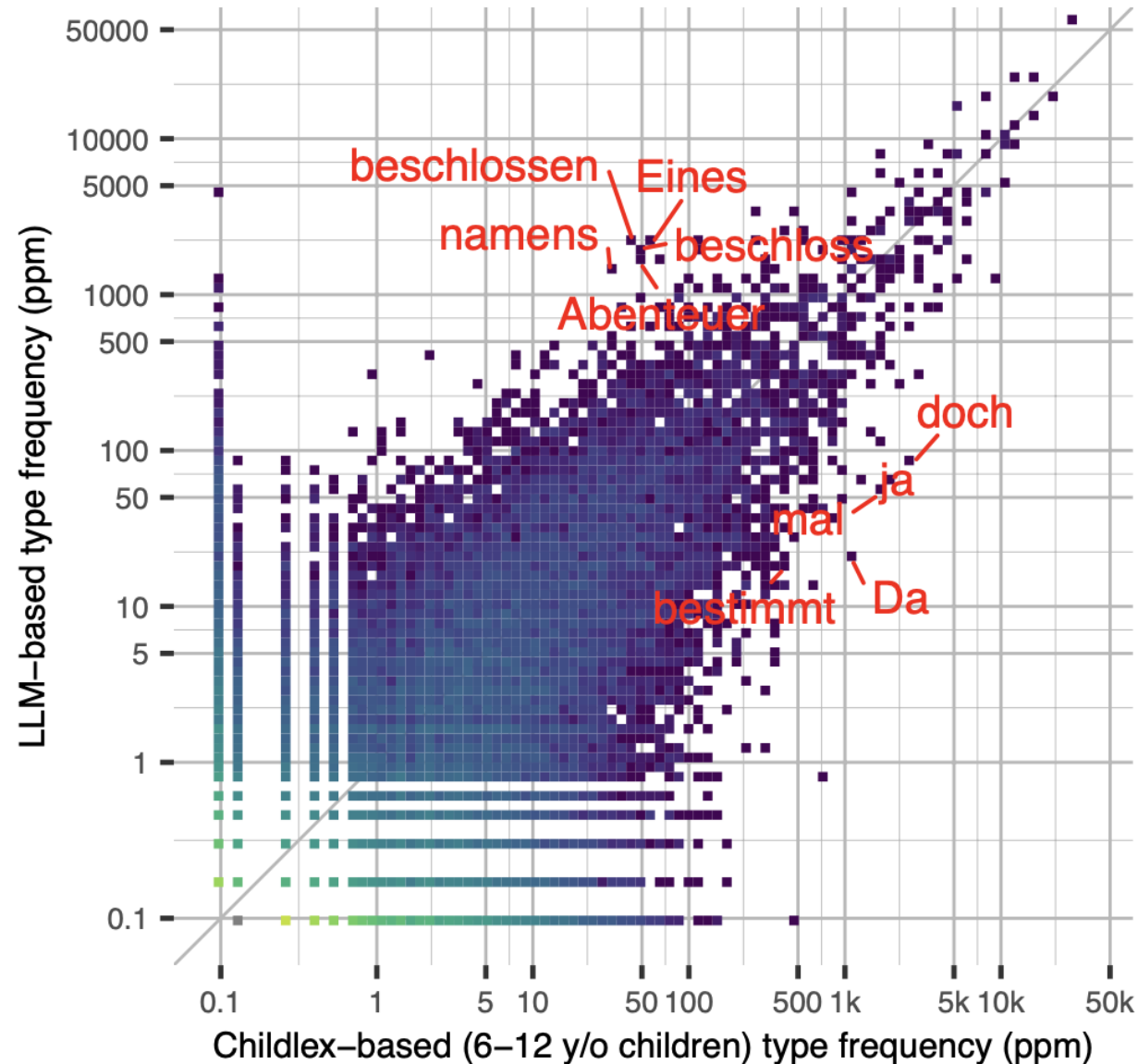
Example text (DeepSeek V3 – long-form text prompt – “ending”)

Mögen wir niemals vergessen, was er uns gelehrt hat: dass die Welt voller Geheimnisse ist, die es zu entdecken gilt, und dass die wahre Stärke oft in den unerwartetsten Gestalten zu finden ist. King-Kong lebt weiter, nicht nur auf der Leinwand, sondern in unseren Herzen und in den Geschichten, die wir von Generation zu Generation weitergeben. Er bleibt ein **Symbol für die ungezähmte Kraft der Natur**, für die Schönheit des Unbekannten und für die unerschütterliche Entschlossenheit, die in jedem von uns schlummert. **King-Kong ist mehr als nur ein Monster oder eine Kreatur** – er ist ein Spiegel unserer eigenen Ängste, Träume und Sehnsüchte.

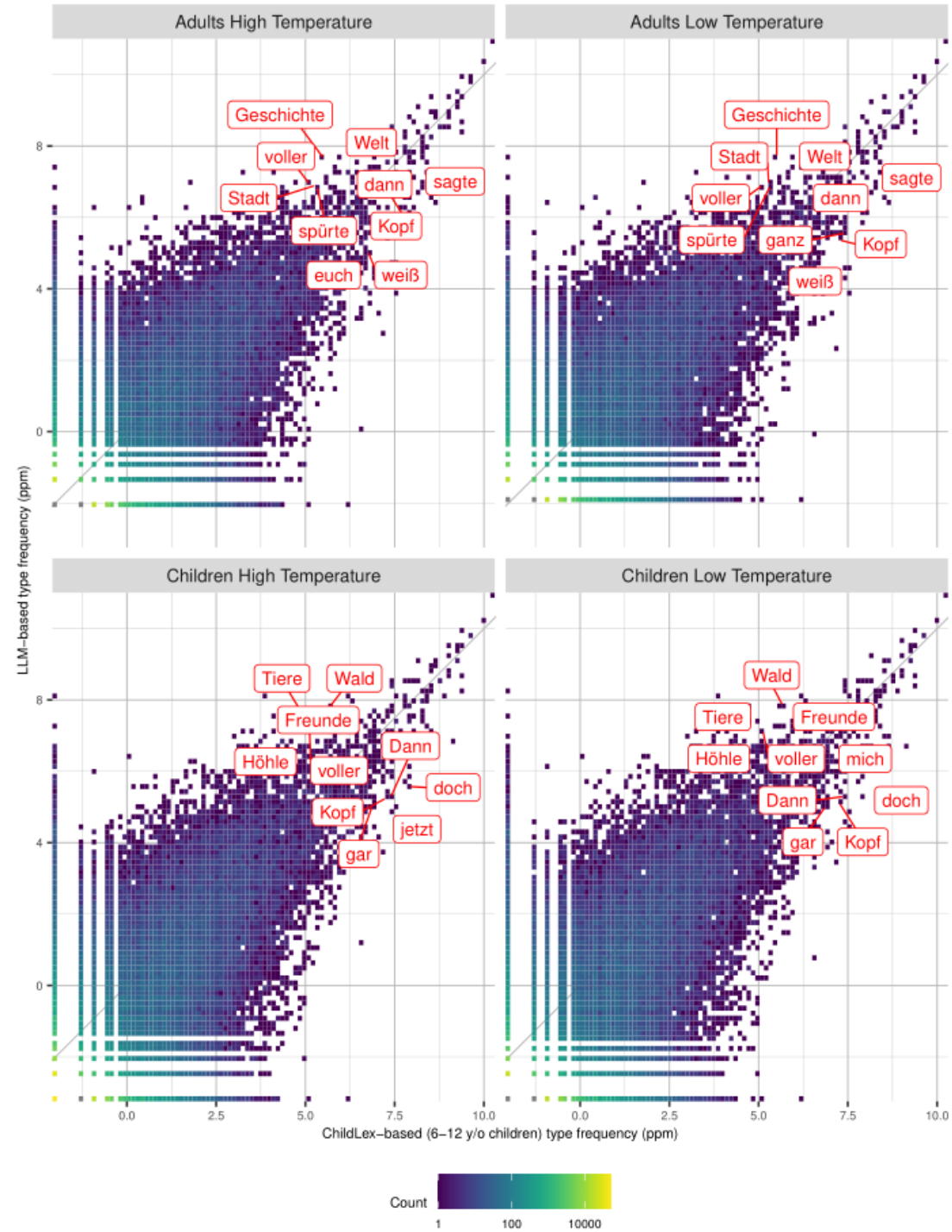
In einer Welt, die oft von Rationalität und Technologie dominiert wird, erinnert er uns daran, dass es noch Raum für das Mystische, das Unerklärliche und das Wunderbare gibt. Er zeigt uns, dass selbst in der Konfrontation mit dem Unbekannten Respekt und Mitgefühl die mächtigsten Werkzeuge sind.

Mögen wir King-Kongs Erbe ehren, indem wir mutig in die unbekannten Welten unserer eigenen Leben vordringen, die Geheimnisse der Natur schützen und niemals aufhören, an die Magie zu glauben, die in jedem Winkel dieser Erde verborgen liegt. Denn solange wir uns an ihn erinnern, wird King-Kong immer bei uns sein – ein stummer Wächter, ein Freund in der Dunkelheit und ein Zeichen dafür, dass die größten Abenteuer oft dort beginnen, wo wir es am wenigsten erwarten.

Word frequency comparison (childLex vs. GPT 3.5)



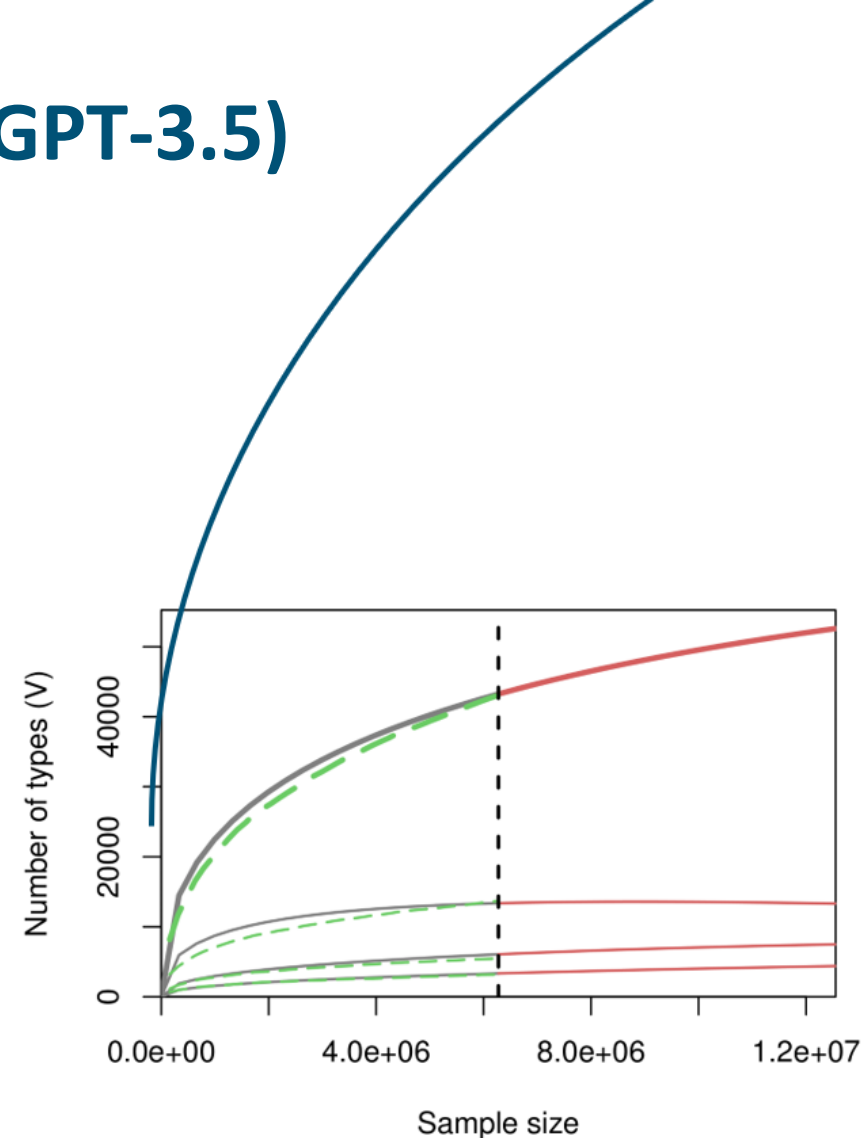
$r = .88$



Lexical richness comparison (childLex vs. GPT-3.5)

Measure	childLex	LLM-corpus
n Books	500	500
Tokens	9,850,786	6,252,808
Types	182,454	46,409
Lemmas	117,952	34,519

Low lexical richness of LLM corpus



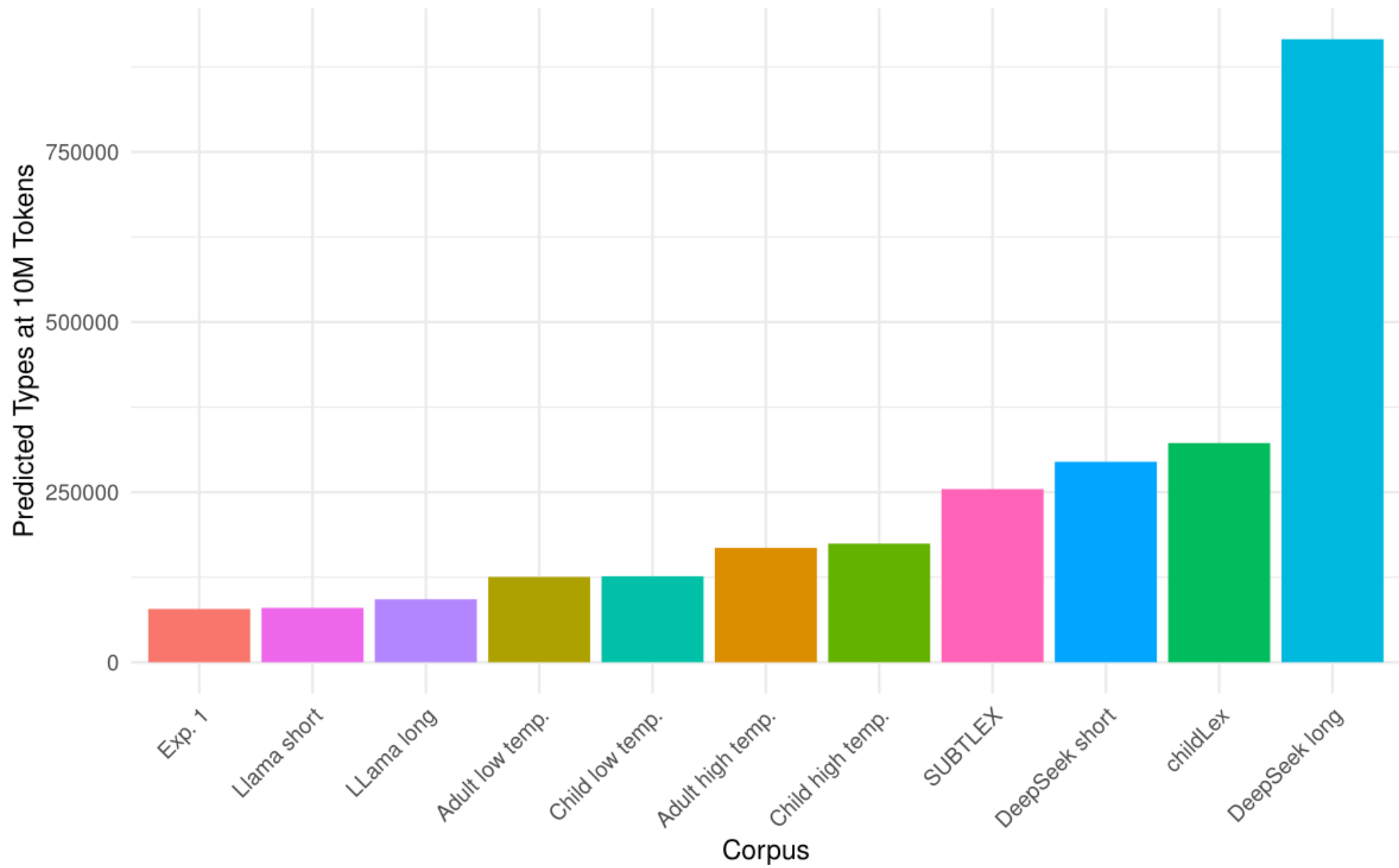
All corpora: Corpus comparison

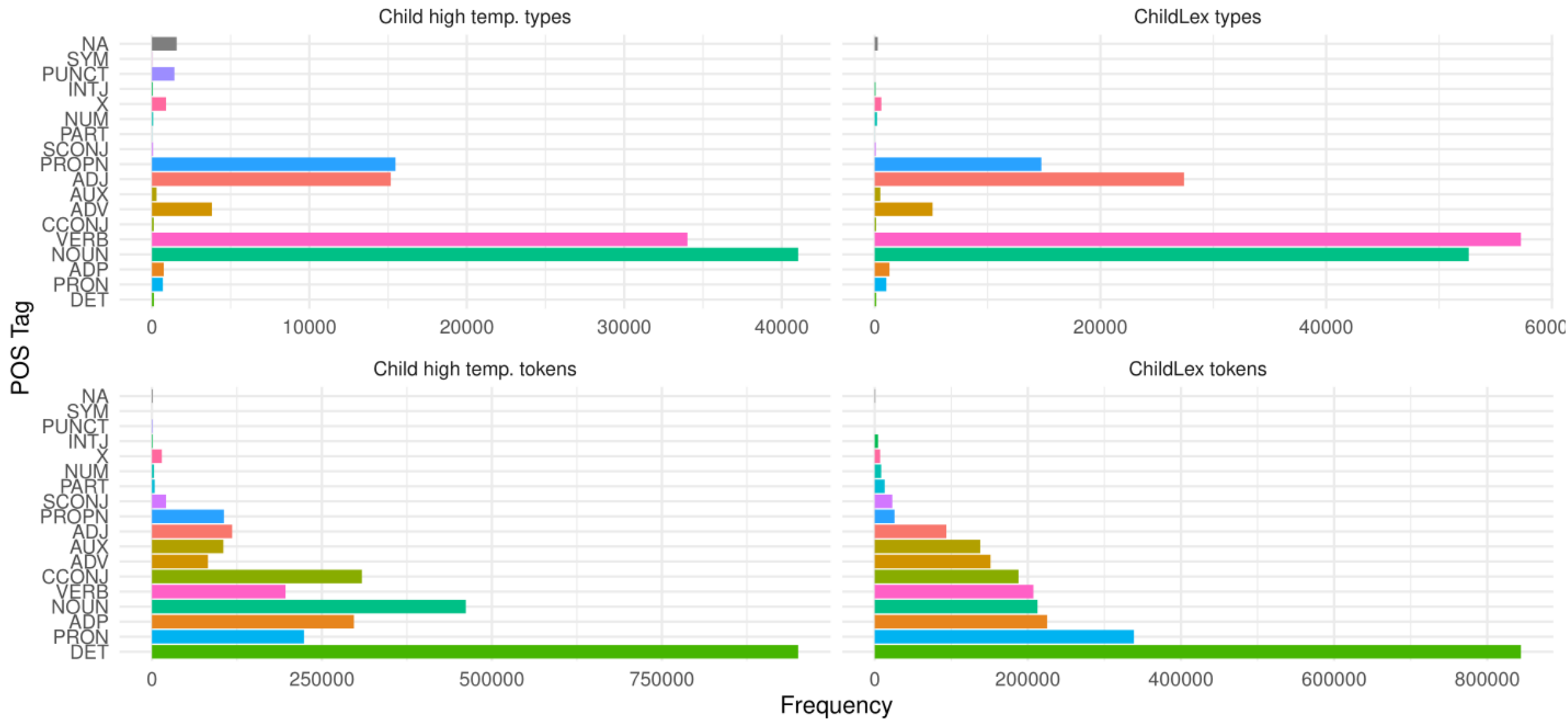
Measure	childLex	LLM-corpus
n Books	500	500
Tokens	9,850,786	6,252,808
Types	182,454	46,409
Lemmas	117,952	34,519

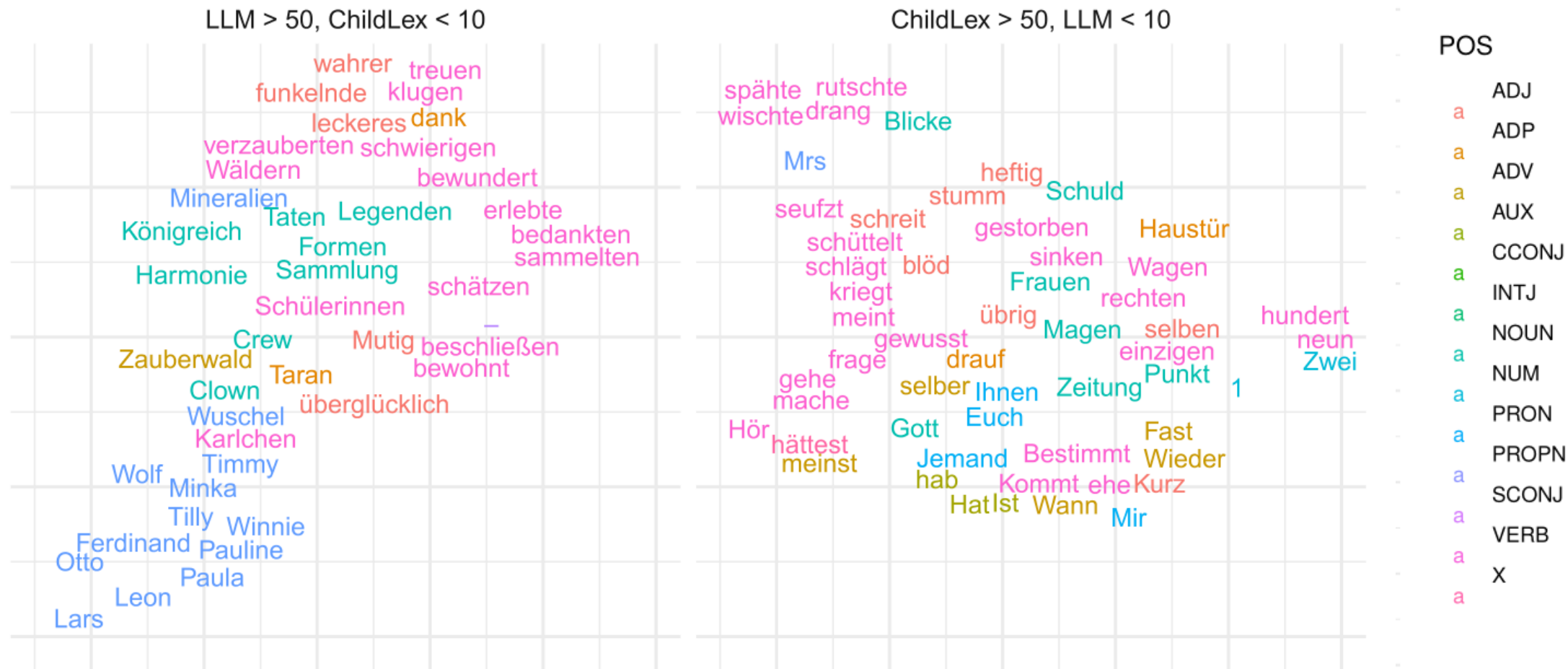
Measure	childLex	Adult Corpus		Child Corpus	
		Low Temp.	High Temp.	Low Temp	High Temp
n Books	500	500	500	500	500
Tokens	9,850,786	7,191,531	7,368,921	23,320,466	23,887,118
Types	182,454	71,423	83,921	84,978	110,603
Lemmas	117,952	52,528	61,318	63,552	82,126

Measure	childLex	Llama long	Llama short	DS-V3 long	DS-V3 short
n Books	500	500	500	500	500
Tokens	9,850,786	10,332,850	7,215,565	9,763,062	7,162,685
Types	182,454	51,320	40,660	239,830	95,321
Lemmas	117,952	39,272	39,272	191,309	74,695

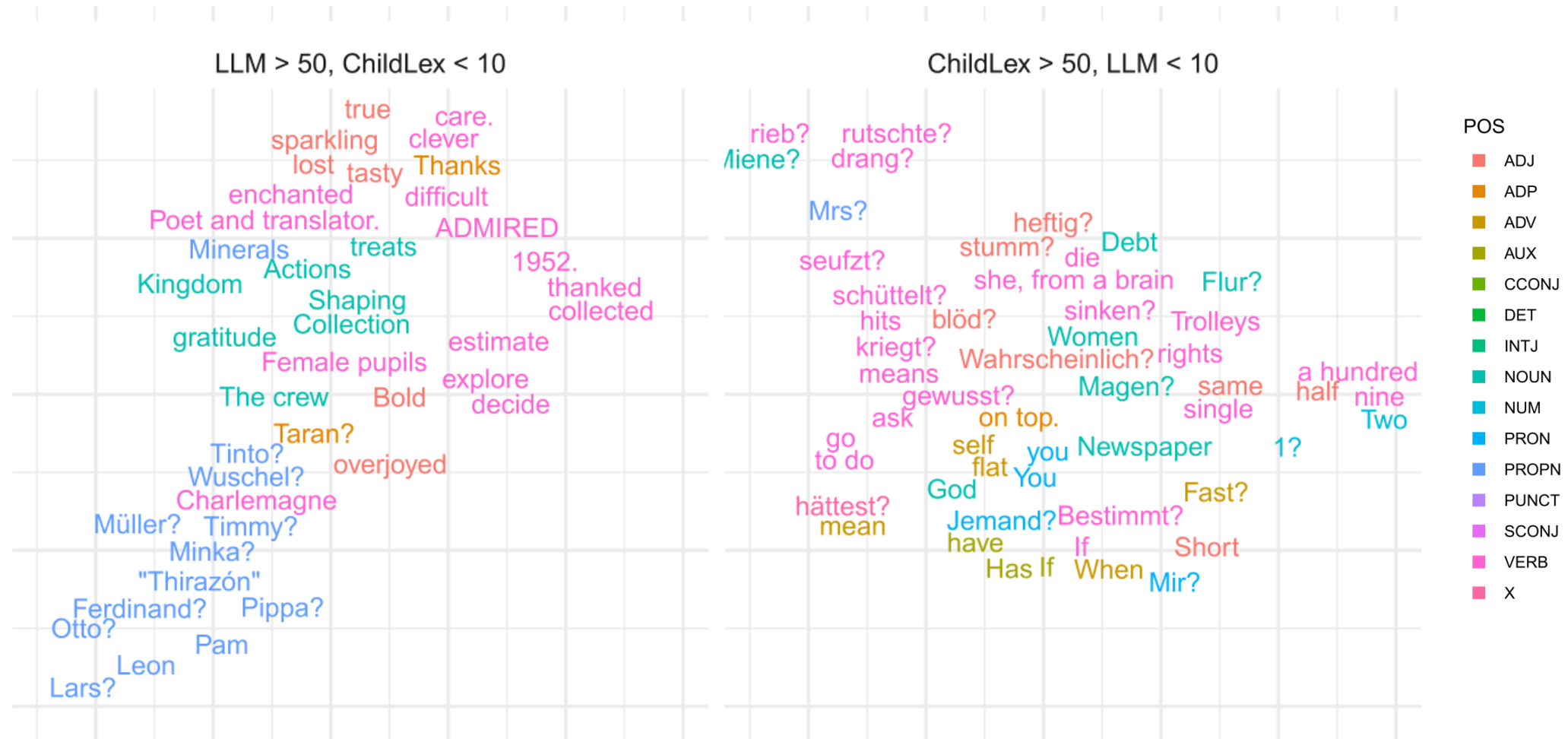
Increase in lexical richness with higher temperature and DeepSeek



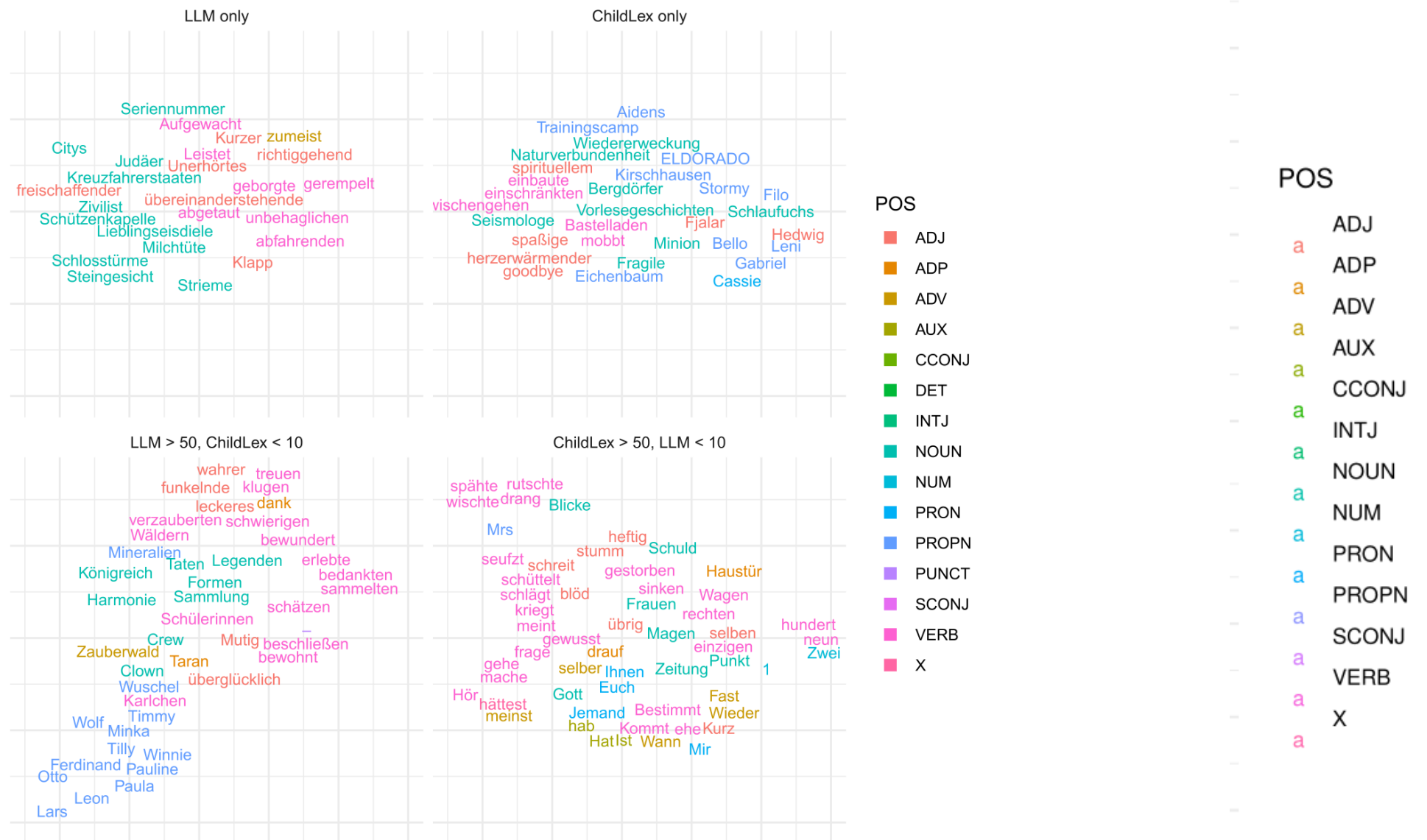




Distributions of word embeddings using dimensionality reduction with UMAP to word vectors from FastText-German. Embeddings were reduced from 300-dimensional vectors. Color coding corresponds to different parts-of-speech (POS) tags



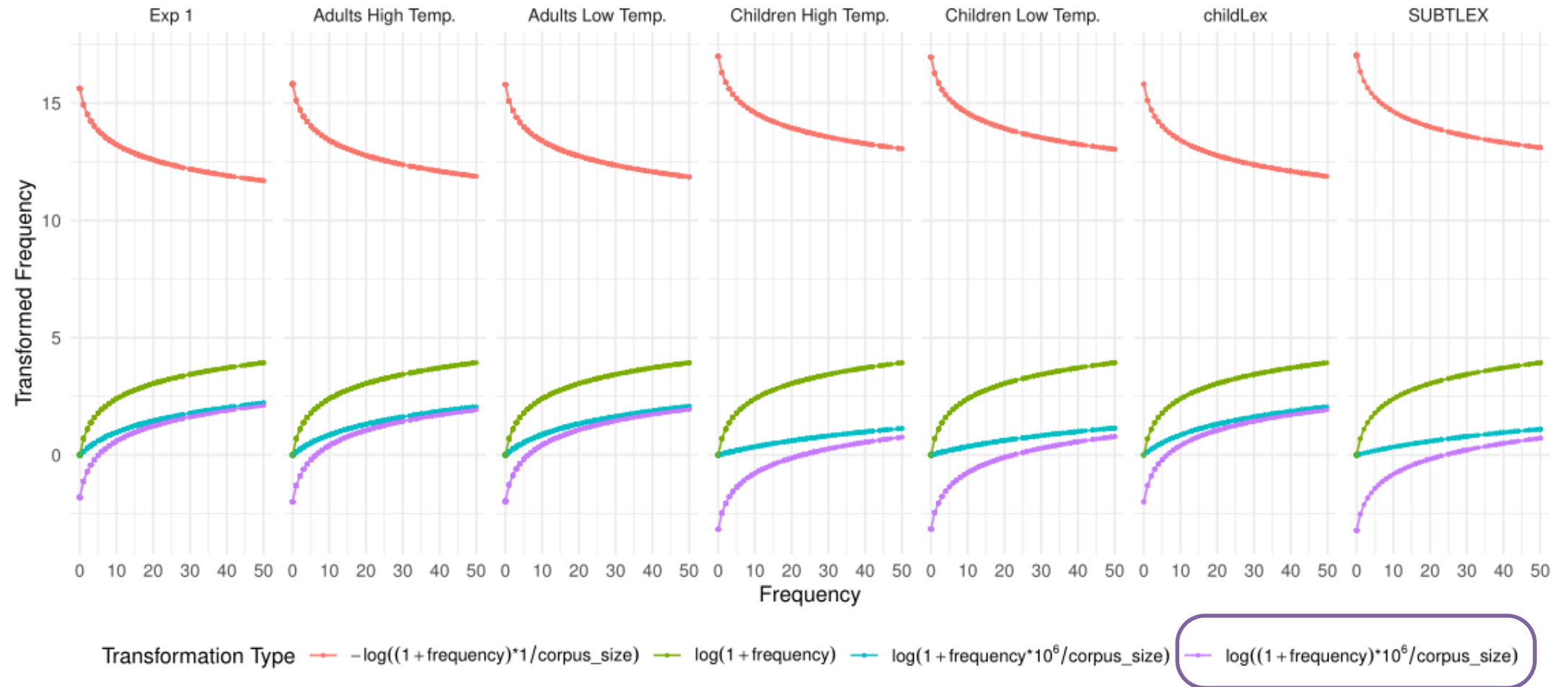
Distributions of word embeddings using dimensionality reduction with UMAP to word vectors from FastText-German. Embeddings were reduced from 300-dimensional vectors. Color coding corresponds to different parts-of-speech (POS) tags



Distributions of word embeddings using dimensionality reduction with UMAP to word vectors from FastText-German. Embeddings were reduced from 300-dimensional vectors. Color coding corresponds to different parts-of-speech (POS) tags

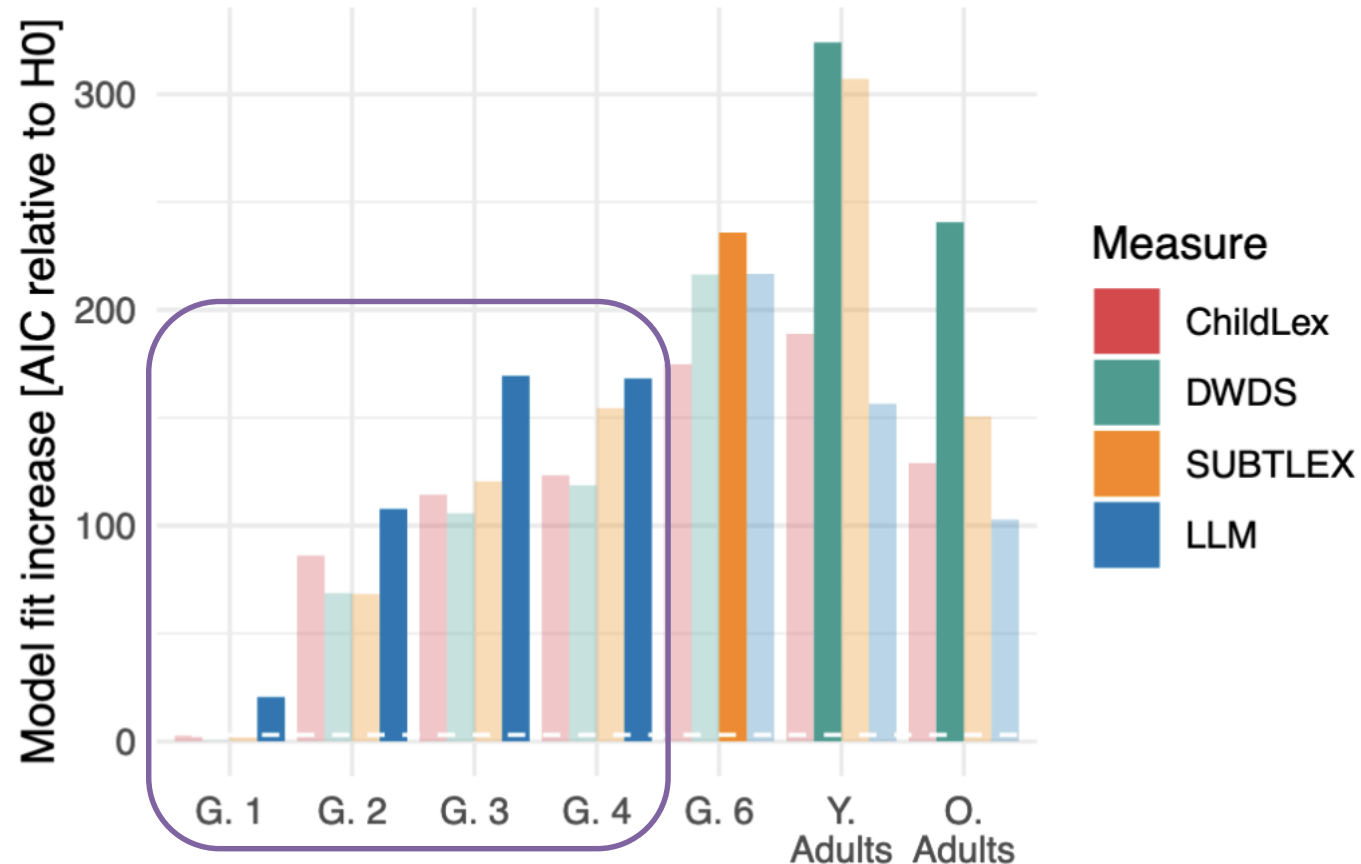
Setup

- Log-transformed normalized word frequency and log-transformed child RTs $\log\left(\frac{(1 + \text{frequency}) \times 10^6}{\text{corpus_size}}\right)$



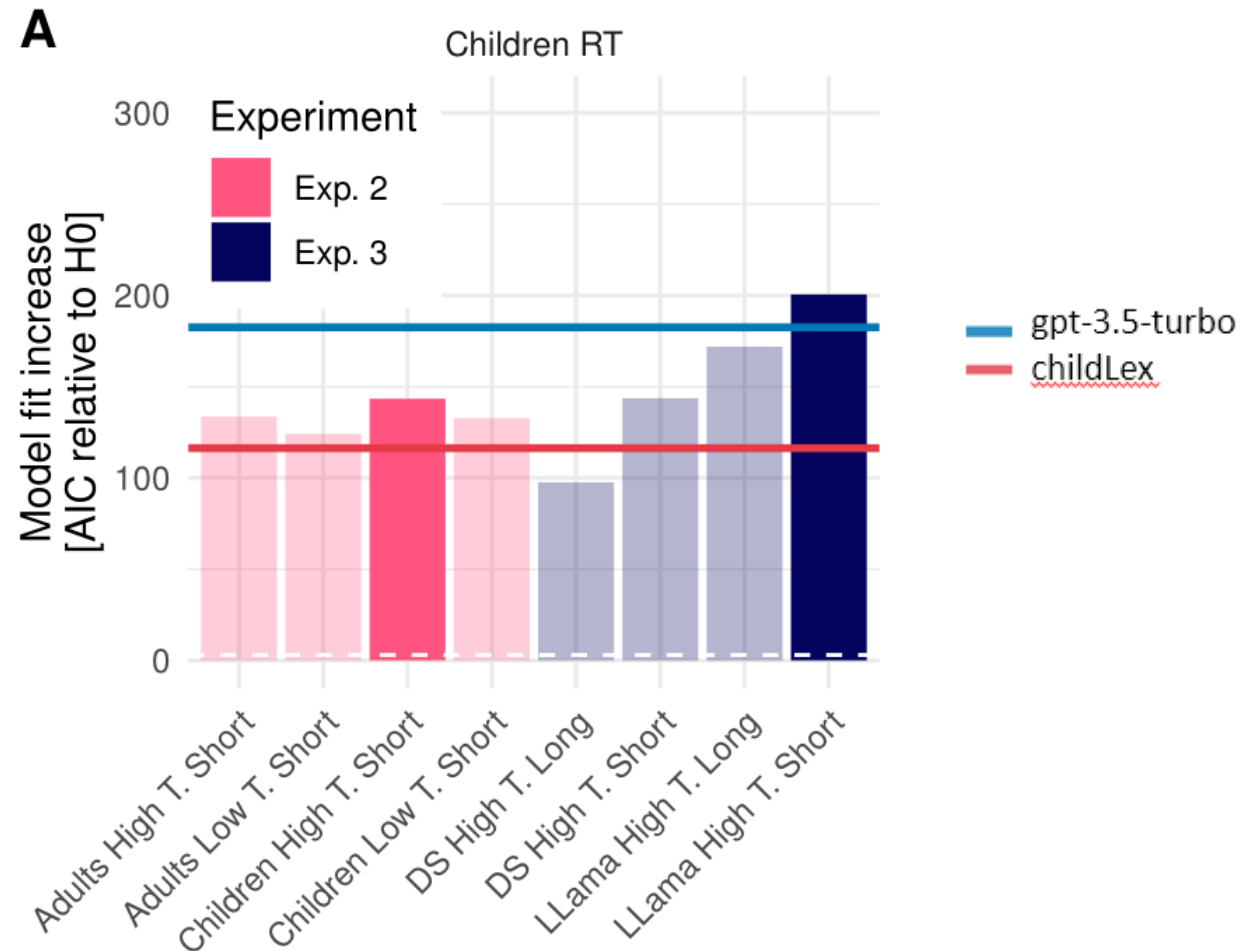
Model fit comparison (childLex vs. GPT-3.5)

RT ~ old20 + aoa + letter.cnt + word.frequency



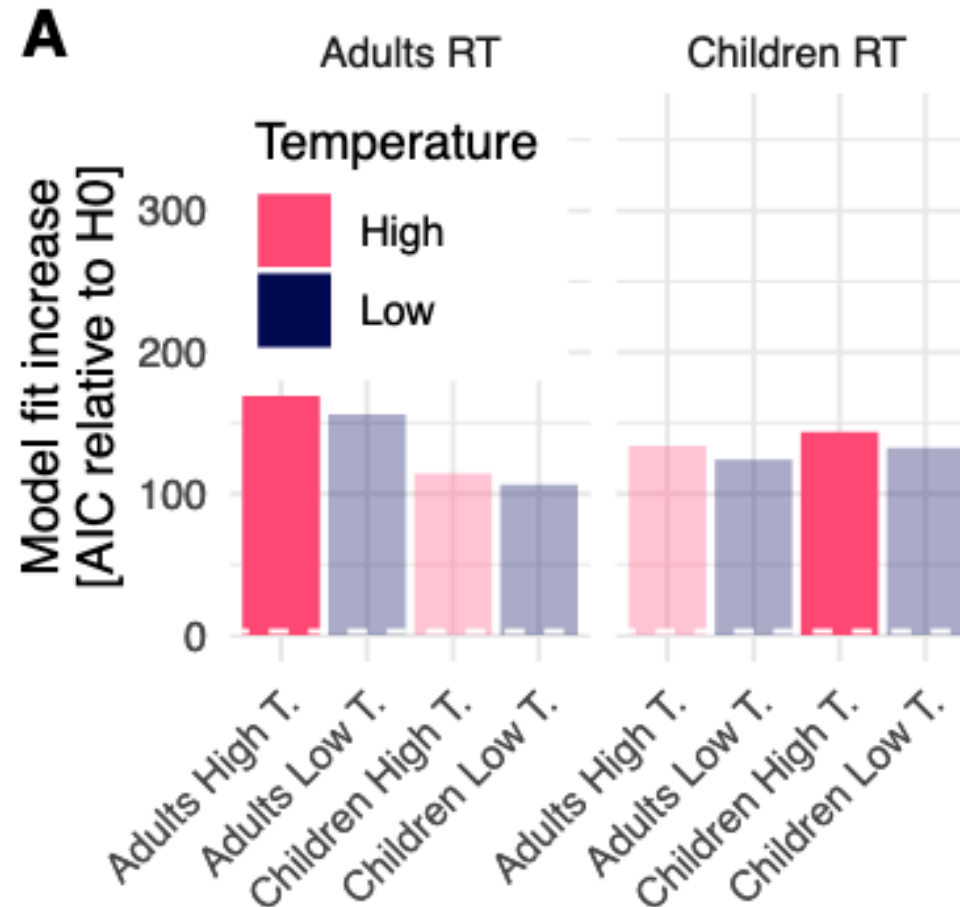
Model fit is higher

All corpora: Estimation of the word frequency effect



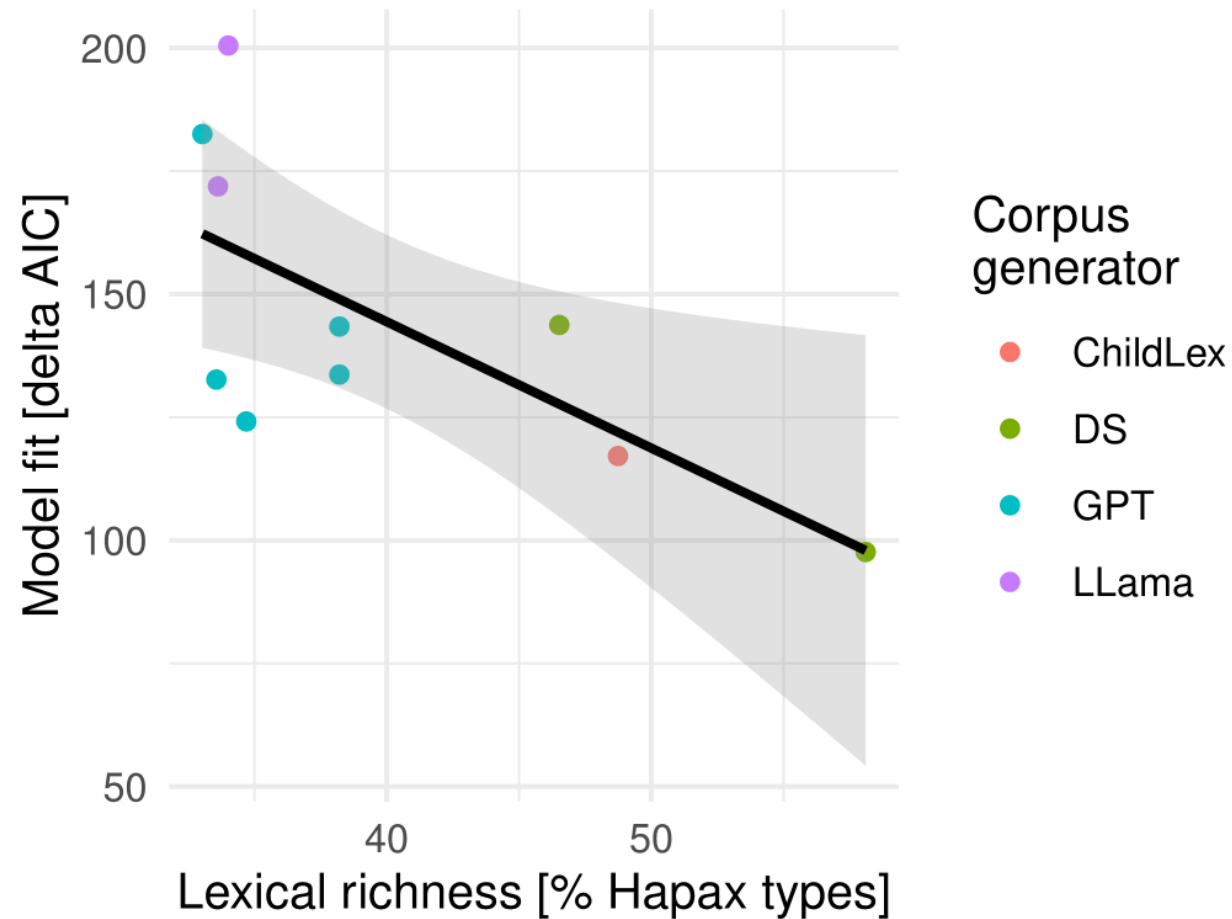
Highest model fit least lexically rich model

All corpora: Estimation of the word frequency effect



Highest model fit for age specific LLM-based frequency

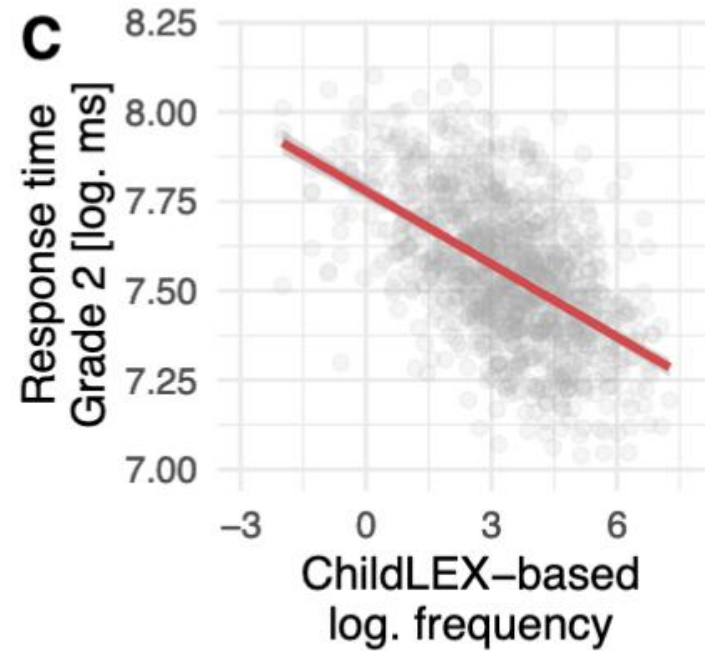
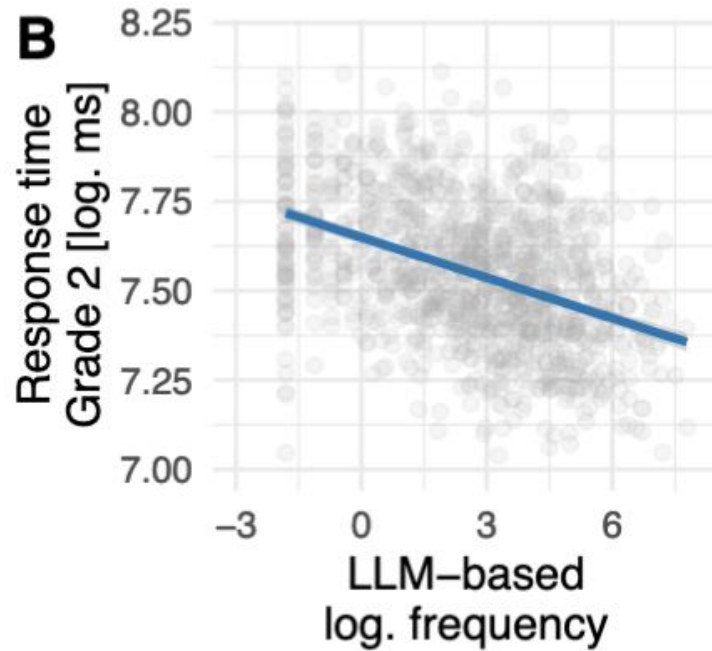
Inverse scaling effect (all corpora)



Less rich → better fit

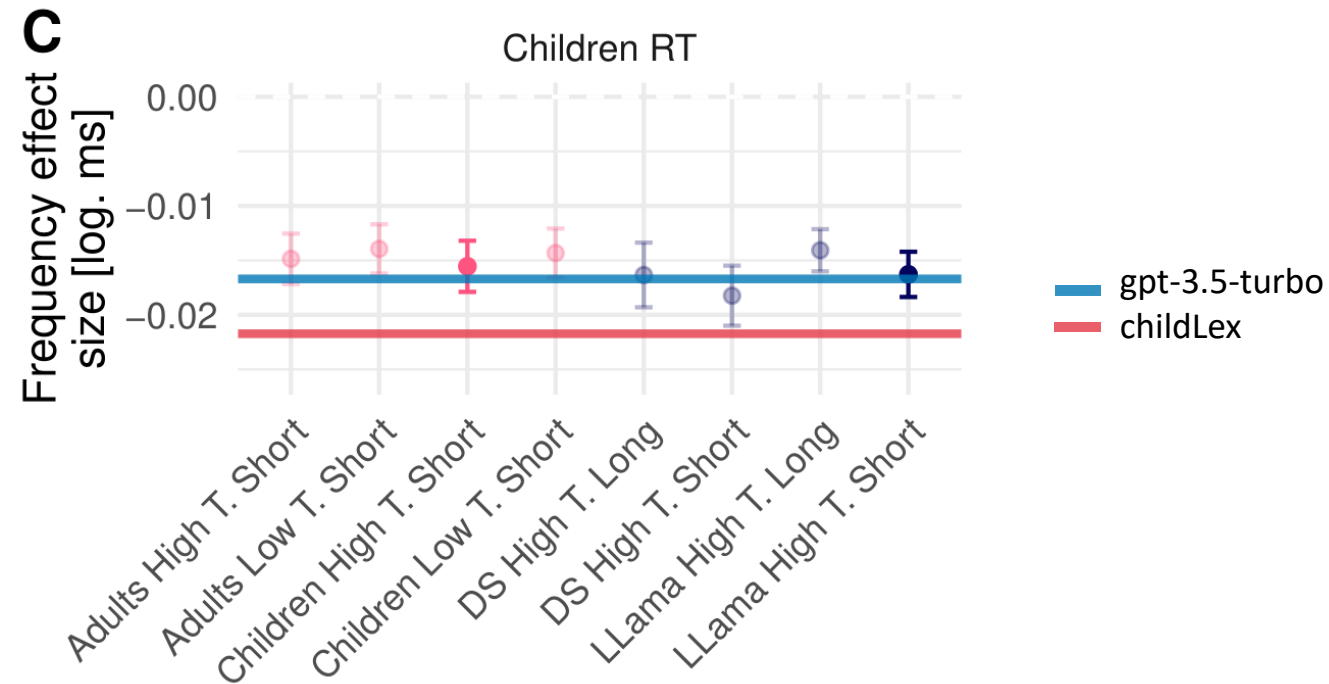
Effect size comparison (childLex vs. GPT-3.5)

RT ~ old20 + aoa + letter.cnt + unigram
+ bigram + trigram + Word Frequency



Effect size is lower

Effect size comparison (all corpora)

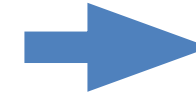


Effect sizes are similar across LLM measures

Summary: Using LLMs for building linguistic corpora

- High **correlation** with childLex word frequency, despite lower richness
 - Better **model fit**, but smaller effect size
 - Temperature & target audience: as expected
 - **Inverse scaling**: Less richness results in better model fit
-
- *Better* representation of word frequency than authors of kids' books?
 - Surprising differences in language use

Find Preprint here



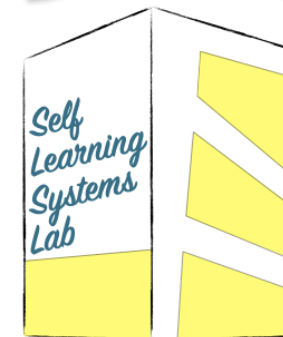
Can large language models generate useful linguistic corpora? A case study of the word frequency effect in young German readers

Job Schepens - Institute for Digital Humanities

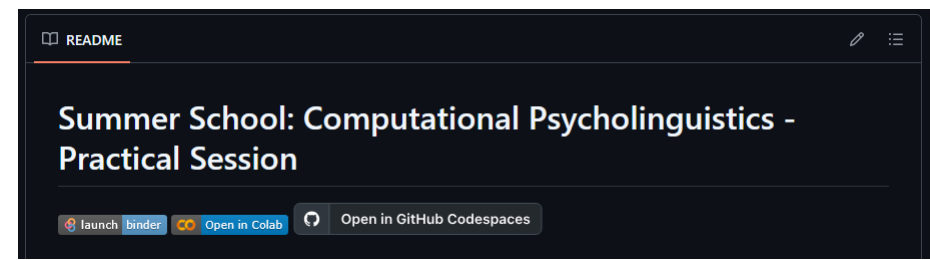
Hanna Wołoszyn - Self learning systems lab

Nicole Marx - Mercator Institute for Literacy and Language Education

Benjamin Gagl - Self learning systems lab

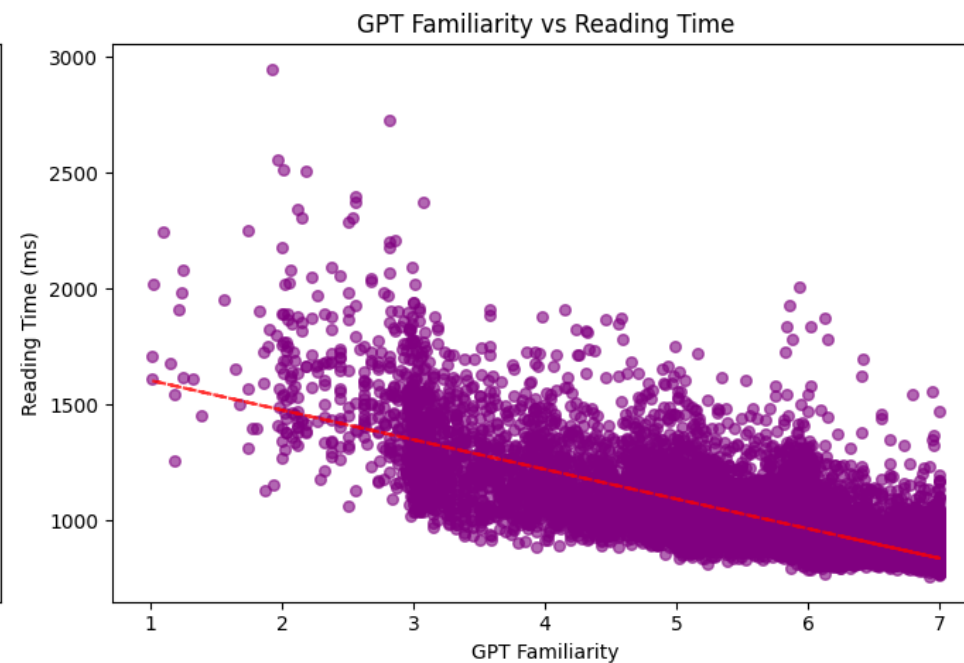
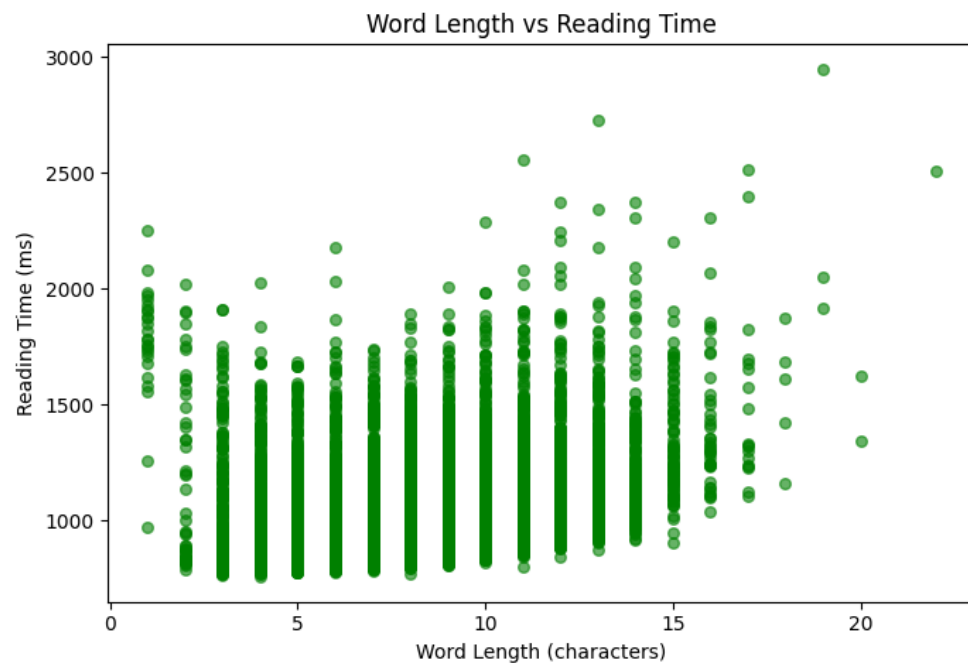
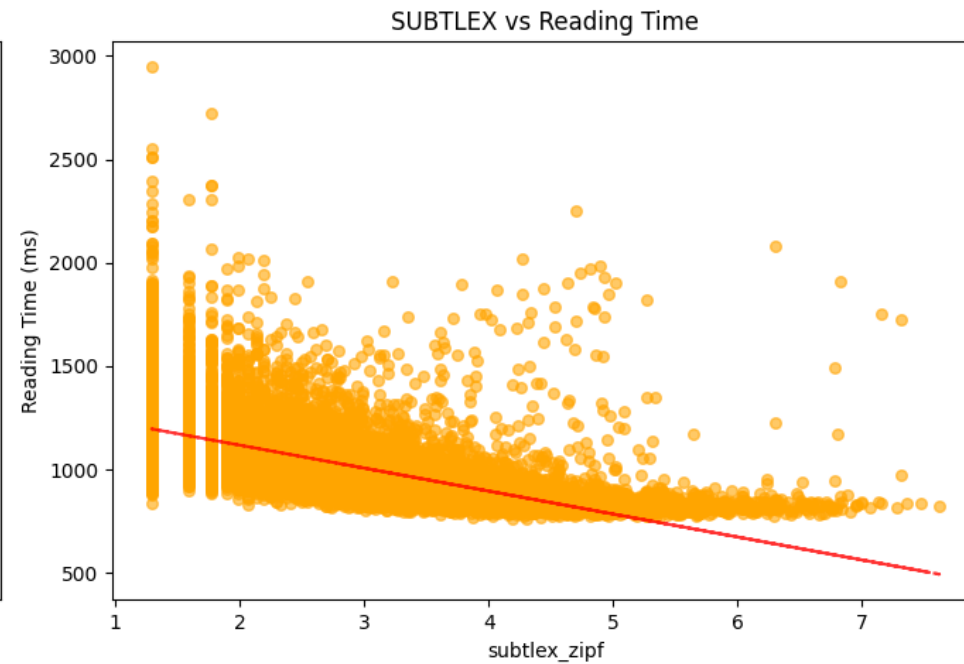
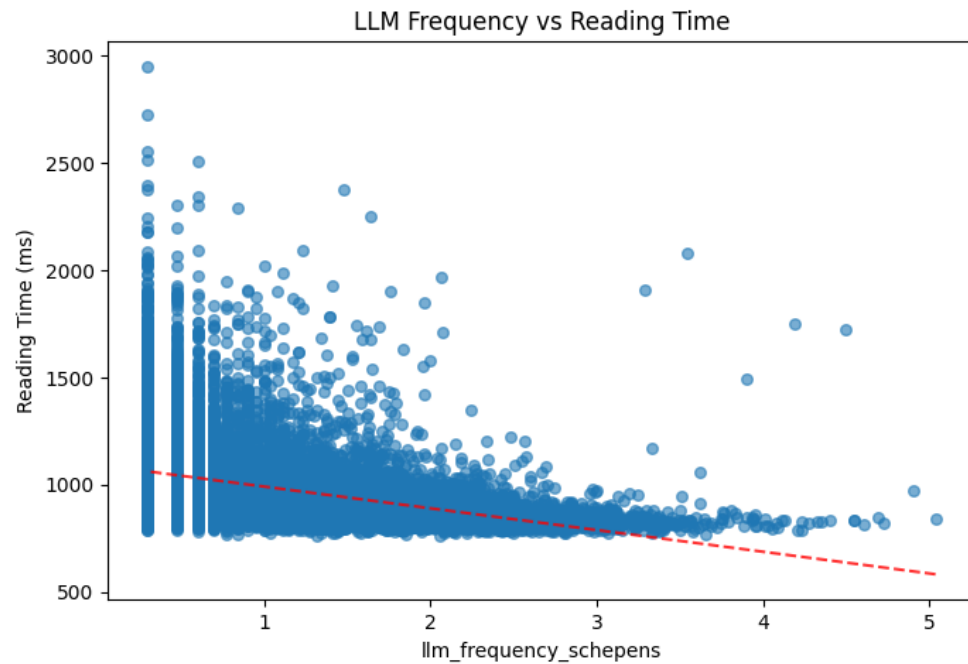


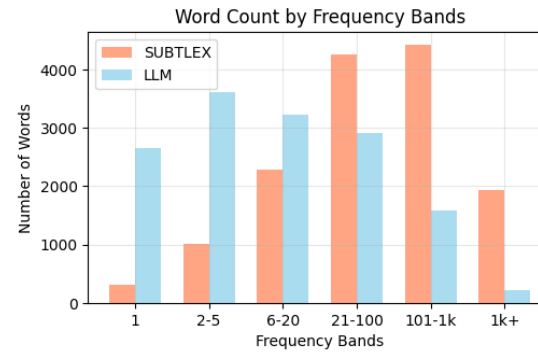
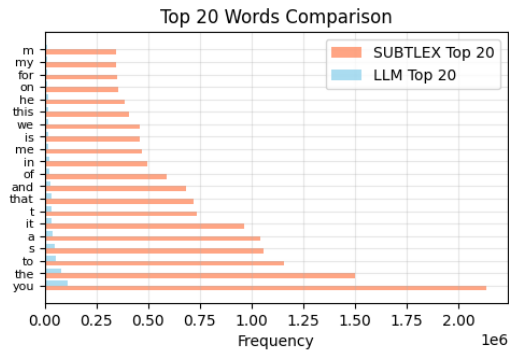
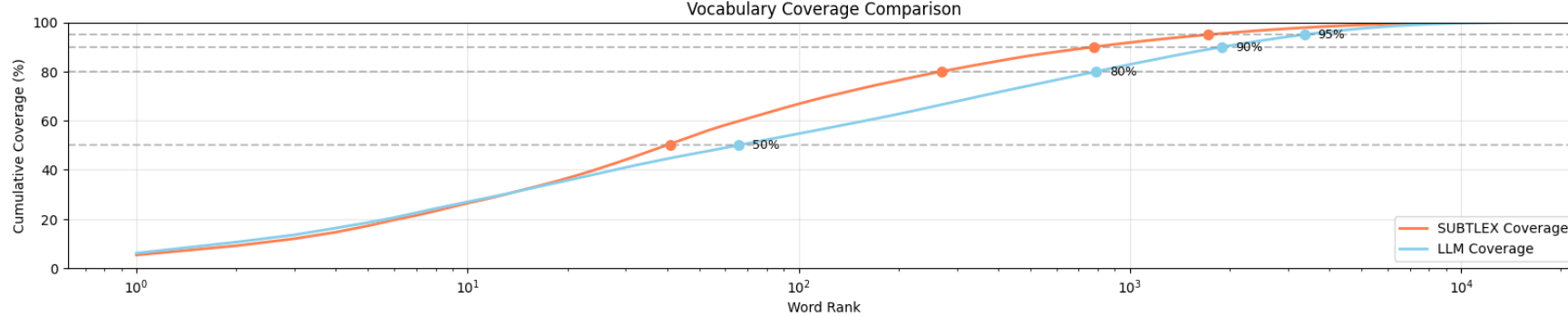
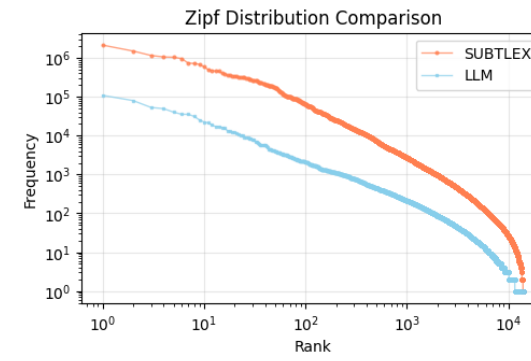
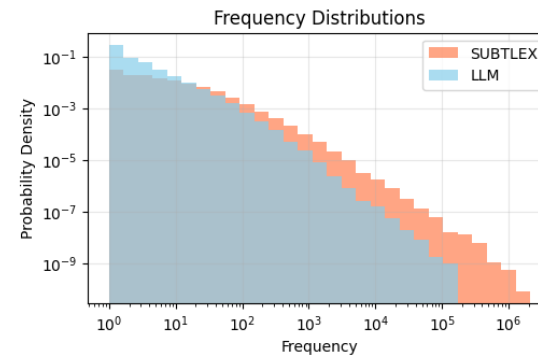
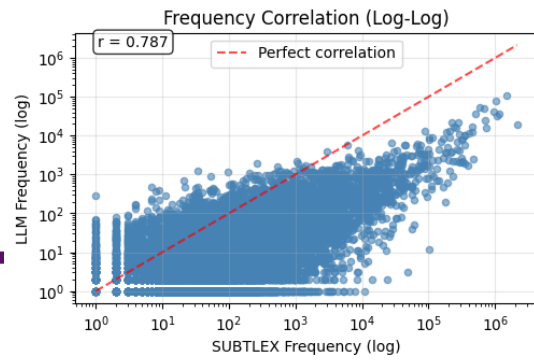
Practice Session: Just run some code?



- Task:
 - **1st part: Notebook 1: Experimenting with LLM-based corpus generation**
 - The repo already contains a few pre-generated 2m corpora (see /scripts folder in the repo)
 - Optional notebooks: Checking the corpus, extracting and formatting with metadata, merging data with behavioral data, comparing different frequency measures, comparing transformations
 - **2nd part: Notebook 2: Validating predictors against human reading times**
- Github repo consists of code and data: <https://github.com/jobschepens/mlschool-text>
 - 40.000 English-speaking adult word reading times
 - Reference data: SUBTLEX, Multilex, GPT familiarity estimates
- Learning goals: try out the pipeline, explore, run new experiments, possibly extend the analysis
 - Understanding how LLMs can be used in computational corpus / psycholinguistics
 - Hands-on: Experiencing pipeline from text generation to statistical modeling
 - Do LLM-based frequencies predict reading behavior better than other frequency measures / familiarity?
- You can run the code locally (recommended) or online
- You can use an Open Router API key to access cheap and fast models such as qwen-30b

LLM Frequency vs Reading Time: Exploratory Analysis





DISTRIBUTION COMPARISON SUMMARY

Sample Size: 14,229 words

SUBTLEX (subtlex_freq_raw):

- Mean: 2815.87
- Median: 76.00
- Max: 2134713
- Min: 1
- Non-zero: 14,229

LLM (llm_frequency_raw):

- Mean: 125.69
- Median: 8.00
- Max: 108943
- Min: 1
- Non-zero: 14,229

