

بهشی پراکتیک 1: ئانالیزی کۆرپس

خولی زمانهوانیی کۆمپیوتهری

<https://sinaahmadi.github.io/KurdishCL/>

سینا ئهحمهدی

ئهركه پراكتيکهکانی ئەم بهشه له بهشی دووی خولهکهدا باس کراون. پيش له ههموو شتيك، کۆرپسيك كه به تايبهت بۆ ئەم خوله کۆ کراوهتهوه و نمونهیهکی بچووکه ليرهدا دابهزینه: بهسته. ئەم کۆرپسه پينج بنزاراوهی کوردی له خۆ دهگرێ: سۆرانی (ckb)، کورمانجی (kmr)، زازا (zza)، ههورامی (hac) و کوردیی باشوور (sdh). له پاش داگرتنی کۆرپسهکان، بهرنامهی AntConc لهم بهستهروهه دابهزینه و کۆرپسيکی تیدا بکهوه.

ئهركهکان

1.0 رینوس و وشه

1. له رینگای کۆرپسهوه، ئەمه بسهلمينه: «هیچ وشهیهکی کوردیی سۆرانی نییه که به ر دهست پێ بکا»
2. له رینگای کۆرپسهوه، ئەمه بسهلمينه: «هیچ وشهیهکی کوردیی سۆرانی نییه که به ل دهست پێ بکا»
3. له خشتهی خوارهوه ههندیك وشه ی هاوواتا هینراون که به شیوازی جیاواز دهنوسرین. له رینگای کۆرپسهوه، فراوانیی وشهکان بدۆزهوه:

فروانی	وشه
	وولات
	ولات
	وهلام
	ولام
	به تايبهتی
	به تايبهتی
	ئهدا
	ئهدات
	دهدا
	دهدات

2.0 دەرپرینه فرەوشەییەکان

دەرپرینه فرەوشەییەکان (multiword expression) بەو وشانەی دەوترێت کە لە چەند بەش پیک هاتوون وەک «ئالوگۆر»، «هەلیت و پەلیت»، «هینانە دەر» و «لە کیسی خۆ دان». بۆ دۆزینەوهی دەرپرینه فرەوشەییەکان دەتوانین لە ئانالیزی کۆرپسەوه کەلک وەرگیرین. هەندیک لە دەرپرینهکان بە یەک شیۆ دەردەکەون، وەک ئەوانەی کە بە شیۆی ناو و کرداریک (ناو + کردار) دەردەکەون وەک «کارکردن» یان «توورە بوون». نموونەیک بەربلاوی کوردی بۆ وشەسازی، کەلک وەرگرتن لە کرداریکی سووک لەگەڵ ناویکە؛ کرداری سووک (light verb) کرداریکە کە لە دەرپرینهکەدا واتاکە بەتەنیا زۆر گرینگ نییە و لەگەڵ وشەکانی دیکەدا مانا دەگری. هەندیک لە کردارە سووکەکان لە سۆرانیدا ئەمانەن:

- «کردن» وەک «بەش کردن» و «پاکردن»
- «دان» وەک «دادان» و «هان دان»
- «بوون» وەک «سوور بوون» و «پابوون»
- «گرتن» وەک «هەلگرتن» و «ماسی گرتن»
- «کەوتن» وەک «سەرکەوتن» و «پەک کەوتن»
- «چوون» وەک «دەرچوون» و «لە دەست چوون»

1. پیتی پیشدانراو (preposition) وەک «لە، بە، دە، وە» و پاشدانراو (postposition) وەک «دا، ڕا، ەوه» بە چ شیوازیک دەردەکەون؟
2. «هوه» وەک لە کرداری «کردنەوه» یان «هاتنەوه» چ رۆلێکی هەیە بە پیتی ئامارەکانی کۆرپسەکه؟
3. دەرپرینه فرەوشەییەکان کە بە کرداری سووکەوه ساز دەکرین بدۆزەوه.

3.0 جیاوازیی نیوان بنزاراوهکان (پیشکەوتوو)

1. 20 وشە پرفراوانی کە لە سۆرانی، کورمانجی، هه‌ورامی و زازادان چینی؟ چۆنیان بەراوەرد دەکەن؟
2. قەبارەیی مامناوەندیی وشەکان لە سۆرانی، کورمانجی، هه‌ورامی و زازا چۆن هەڵدەسەنگینی؟
3. وشەیی هاوڕێ (collocation) بەو وشانەی دەوترێت کە بە فراوانییەکی زۆرەوه پیکەوه دین و مانایەکی جیاوازیان هەیە. بۆ نموونە، ئەمانە وشەیی هاوڕێن: «چەپکیک گۆل» نەک «کۆمەلێک گۆل»، «ئاو خوارنەوه» نەک «ئاو خواردن» و «وەرینی سەگ» نەک «وەرینی بەرخ». بەپێی کۆرپسەکه، هەندیک لە وشە هاوڕێیەکانی بنزاراوهی کۆرپسەکه بدۆزەوه.
4. یاسای زیپف (Zipf's law) لە ئانالیزی کۆرپسدا دەلی ئەگەر وشەکان لە کۆرپسێکدا بە پیتی فراوانییان ریزبەندی بکەن، وشەیی یەکەم دوو ئەوەندە وشەیی دووهم فراوانی هەیە و سێ ئەوەندە وشەیی سێهەم و هتد. کەوا بوو، ئەو وشەیی کە لە ریزی n دایە، $\frac{1}{n}$ جاری یەکەم وشە فراوانی هەیە. دەتوانی ئەمە لە کۆرپسە سۆرانی، کورمانجی، هه‌ورامی و زازاکەدا بسەلمینی؟

- [1] Sina Ahmadi. Building a corpus for the Zaza–gorani language family. In *Proceedings of the 7th Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 70–78, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics (ICCL). <https://aclanthology.org/2020.vardial-1.7.pdf>.
- [2] Sina Ahmadi. A tokenization system for the Kurdish language. In *Proceedings of the 7th Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 114–127, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics (ICCL). <https://aclanthology.org/2020.vardial-1.11.pdf>.
- [3] Heather Froehlich. Corpus analysis with antconc. *Programming Historian*, 2015. <https://programminghistorian.org/en/lessons/corpus-analysis-with-antconc>.