

## KORPUSDAN QIDIRUV TIZIMI SIFATIDA FOYDALANISH

**O‘roqova Dilnoza Odil qizi**

Mirzo Ulug‘bek nomidagi O‘zbekiston  
Milliy universiteti 2-bosqich magistranti.

E-mail: [adilovnadilnoza09@gmail.com](mailto:adilovnadilnoza09@gmail.com)

**Annotatsiya:** Mazkur maqola korpus tushunchasi, uning tildagi o‘rni, korpusning qidiruv turlari, ishlash prinsipi va ularni qidirish, o‘rganishga bo‘lgan turlicha yondashuvlarni o‘z ichiga oladi.

**Kalit so‘zlar:** korpus, kontekst, lug‘at, intertekstuallik, yashirin modellar, tillarni taqqoslash, metama‘lumot, matn razmetkasi va lingvistik annotatsiya

**Abstract:** This article includes the concept of the corpus, its place in the language, types of corpus searches, the principle of operation and different approaches to their search and study.

**Keywords:** corpus, context, vocabulary, intertextuality, implicit models, language comparison, metadata, text markup and linguistic annotation

**Аннотация:** В данной статье представлены понятие корпуса, его место в языке, виды корпусных поисков, принцип работы и различные подходы к их поиску и изучению.

**Ключевые слова:** корпус, контекст, словарный запас, интертекстуальность, неявные модели, языковое сравнение, метаданные, текстовая разметка и лингвистическая аннотация.

Dunyo tilshunosligida tilning imkoniyatlarini kengroq o‘rganish, til grammatikasining muammoli tomonlarini kontekstda aniqlash, tilda grammatik qoliplarni belgilash, ko‘p tarmoqli elektron lug‘atlar yaratish ishini yengillashtirish, tilni o‘rganishda zamonaviy axborot texnologiyalardan foydalanish samaradorligini oshirish, tilda avtomatik tarjima, qidiruv va komyuter tahlilini yo‘lga qo‘yish, elektron darsliklar va lug‘atlar tayyorlash kabi masalalarini hal qilish uchun tillarda korpuslar yaratishning nazariy va amaliy asoslarini ishlab chiqish, tilning maxsus sohalar bo‘yicha korpusini qurish zaruratining mavjudligi maqolamizning dolzarbligini belgilaydi. Korpus matnlarning katta to‘plami bo‘lishi bilan birga, bu lingvistik tahlilga asoslanadigan yozma yoki og‘zaki materiallar asosi hamdir. Korpusda:

- 1) aniq bir yozuvchining yoki yozuvchilarning matnlari;
- 2) aniq o‘n yil yoki asrga oid matnlar;
- 3) muayyan mavzudagi zamonaviy matnlar;
- 4) til yoki jamiyatda yetarlicha mavjud bo‘lgan zamonaviy matnlardan foydalaniladi.

Yuqorida belgilangan matn turlaridan korpus shakllantirilganda quyidagilarni qidirish mumkin:

- 1) kontekstdagi so‘zlarning barcha shakllari;
- 2) lug‘atdagi o‘zgarish va izchillik;
- 3) eng ko‘p tanlangan so‘z birikib kela oladigan so‘zlar;
- 4) ikkita matnlar o‘rtasidagi eng muhim farqlar;

- 5) muayyan yozuvchining soʻz va iboralami qoʻllashdagi oʻziga xoslik;
- 6) intertekstuallik: soʻzning maʼnosi uning qoʻllanilish yigʻindisi sifatida;
- 7) soʻz birliklaridan foydalanadigan yashirin modellar;
- 8) tillarni taqqoslash.

Korpusning qidiruv tizimi va uning ishlash prinsipi 3 turdagi asosiy qismni oʻz ichiga oladi, bu korpusdagi maʼlumotlar bazasini tekshirishga yordam beradi. Bular: metamaʼlumot, matn razmetkasi va lingvistik annotatsiya (izohlash). Metamaʼlumotlar - matn kim tomonidan yozilgani, qachon nashr qilingani, qaysi tilda yozilgani toʻgʻrisida baʼzi maʼlumotlarni beradi. Metamaʼlumot korpus matnida kodlanishi yoki alohida hujjat sifatida maʼlumotlar bazasida saqlanishi mumkin. Matn razmetkasi yoki matnni belgilash matnni formatlashni ifodalash uchun ishlatiladi masalan, jumlaning boshlanishi va tugash nuqtasini belgilashda muhim hisoblanadi. Metamaʼlumot matndagi maʼruzachilarni aniqlab beradi va har birining yoshi, jinsi haqida foydali maʼlumot beradi. Matn razmetkasi keyinchalik har bir maʼruzachi gapirishni boshlaganini va tugatganligini koʻrsatish uchun ishlatiladi. Korpusga kiritilgan metamaʼlumot va matn razmetkasi birgalikda bir qator tadqiqot savollariga javob beradi. Matn korpusi ichida lingvistik maʼlumotlarni kodlashimiz mumkin, keyinchalik bu tahlilni muntazam va aniq tarzda olib boriladi, bu holatda korpus analitik yoki lingvistik nuqtayi nazar bilan izohlanadi. Annotatsiya kodlash qoidalaridan matn razmetkasi sifatida foydalanadi, masalan, XML (extensible markup language) qirrali-briket teglari ishlatilishi mumkin, bu yerda otli fraza boshlanadigan nuqtada (<np>) shakli va tugash nuqtasida (</np>) shakli ishlatiladi. Masalan, *np Oʻquvchi </np> <np> doskaga </np> chiqdi (<np> The children </np> sat in <np> the room </np>)*[Abdurakhmonova, N. 2019,2021,2022].

Korpus lingvistikasida korpuslar foydalanuvchiga tez va ishonchli qidirishga imkon beradigan dasturlarni qoʻllashda oʻzgarimlidir. Bu dasturlarning baʼzilari, chunonchi, konkordanslar foydalanuvchiga matndan soʻzlarni qidirishga imkon beradi. Shunga oʻxshash dasturlarning koʻpchiligi matnlardagi maʼlumotlar chastotasini yaratishga imkon beradi. Masalan, korpusda paydo boʻladigan barcha soʻzlarning roʻyxati va har bir soʻzning maxsus oʻsha korpusda necha marta uchrashini koʻrsatadigan soʻzlarning chastotali roʻyxatini yaratadi. Konkordans va chastota faktlari tahlilning ikki muhim tomoni hisoblanadi, yaʼni sifat va miqdoriy jihati sifatida korpus lingvistikasining muhim qismlaridir.

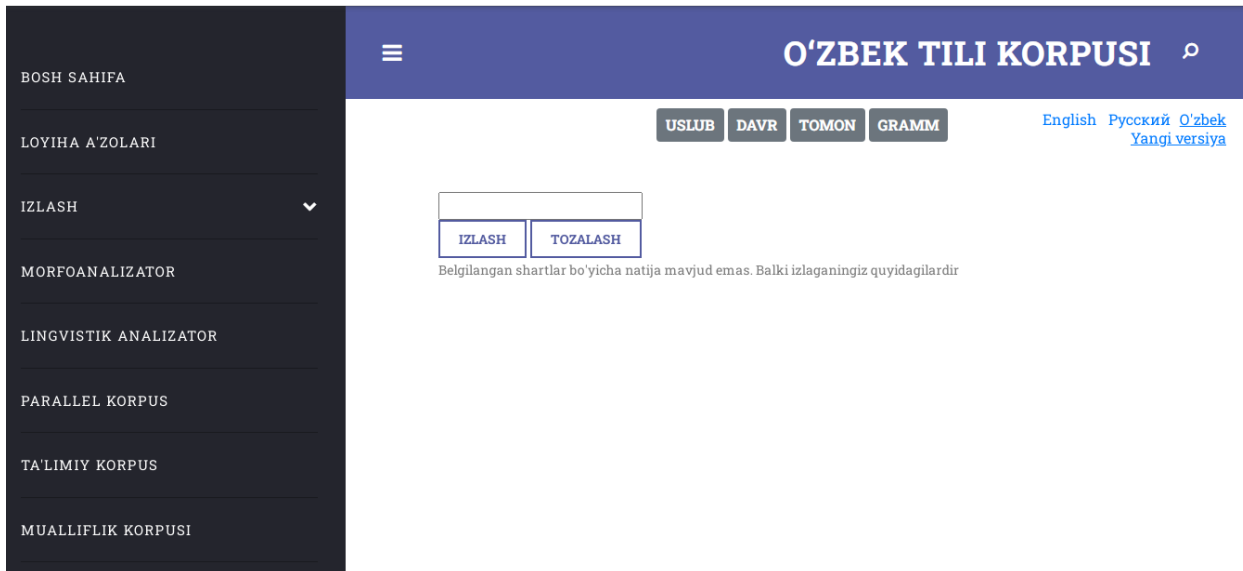
Korpus qidiruv tizimidagi eng muhim saytlardan biri *uzbekcorpus.uz*. Oʻzbek tili korpusi uchun lemme va token asosiy qidiruv birligi sanaladi. Lekin korpusning qidiruv menejeri ham token, ham soʻz birikmasi va gap qidiruviga moʻljallangan. Lingvistik maʼlumotlar bazasida oʻzbek tilining morfologik bazasi boʻlib, barcha qidiruv algaritmi uchun asos vazifasini bajaradi. U orqali korpusda mavjud soʻzlarning statistikasini chiqarish imkoniyati tugʻiladi. Shuningdek, korpusning qidiruv tizimi orqali ixtiyoriy soʻzning n-gram boʻyicha qoʻllanilishi natijasida soʻzlarning konkordanslarni hosil qilish mumkin. Lemma, soʻzning lugʻatdagi shakli boʻlib, aksariyat tillar uchun yasaliq va grammatik jihatdan shakllanish materiali boʻla oladi. Oʻzbek tilining turkiy tillarga xos xususiyatlardan biri morfem birliklarini kombinatsion shakllanishi boʻyicha natijaga erishildi[Abdurakhmonova, N. 2019,2021,2022].

Oʻzbek tili elektron korpusning funksiyalari quyidagidan iborat:

- 1) Bosh sahifa
- 2) Izlash (lemme, token, konkordans boʻyicha)
- 3) Morfoanalizator

- 4) Subkorpuslar (ta'limiy, parallel korpus)
- 5) Foydalanuvchilar yo'riqnomasi
- 6) Elektron lug'atlar
- 7) Tezaurus

Korpusning qidiruv tizimi uch asosiy birlik bo'yicha funksional vazifa bajaradi



- 1) Matnlarni uslub bo'yicha qidirish
  - 2) Matnlarni xronologik ya'ni davr tartibida qidirish
  - 3) Konkordanslar uchun n-gram modeli bo'yicha o'ng va chap tarafdin qidirish
- Ushbu elektron korpus statistic ma'lumotlarni quyidagi jihatlarga ko'ra chiqarish imkoniyatiga ega:

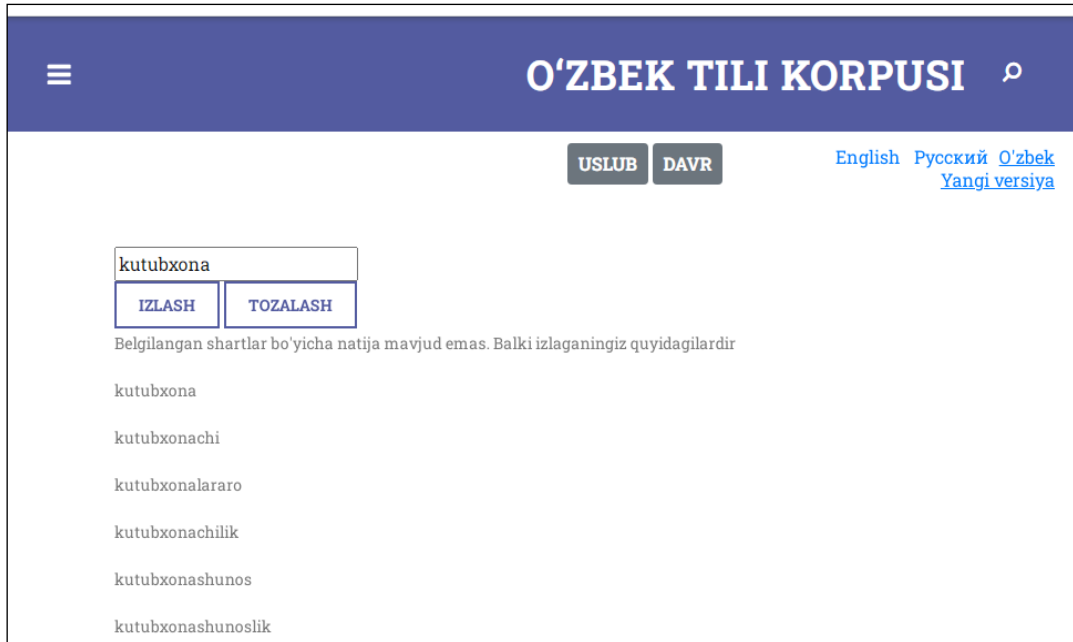
- 1) lemma bo'yicha qidirish – masalan, kitob so'zi qidirilganda, ushbu so'z o'zak holatda matn tarkibida aniqlanadi



Lemma bo'yicha qidirishda so'zlarni lug'at shaklida 0 grammatik ko'rsatkichda yoki aynan grammatik qo'shimchalar bilan qidiruvda so'rovga berilgan shakliga qarab natija uzatiladi [Abdurakhmonova, N. 2019,2021,2022].

Lemma korpusdan ma'lumot qidirishning asosiy birligi sanaladi. Morfologik tahlil natijasida mustaqil so'z turkumining statistik natijasi olinadi. Morfologik shakllangan so'zlar kombinatsion jihatdan turli shakllarga ega bo'ladi.

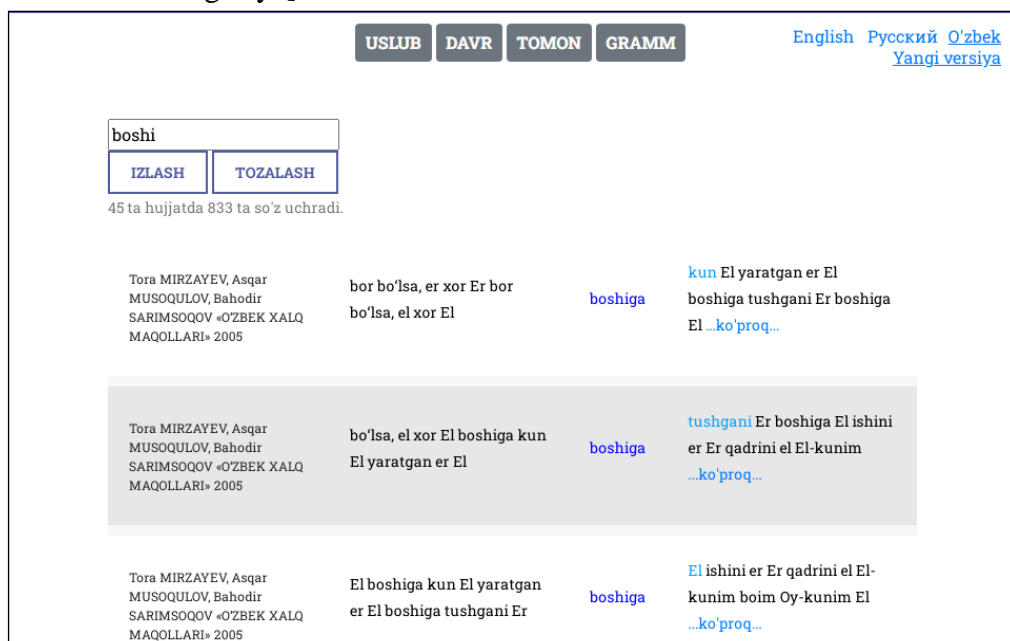
2) token bo'yicha qidirish – [ asos + grammatik kategoriya] modeli bo'yicha so'rov yaratish. Mazkur holatda yasovchi qo'shimchalar yoki shu so'z qatnashgan birlik morfologik jihatdan tahlil qilinadi



The screenshot shows the search interface of the O'zbek Tili Korpusi. At the top, there is a navigation bar with a menu icon, the title "O'ZBEK TILI KORPUSI", and a search icon. Below the navigation bar, there are buttons for "USLUB" and "DAVR", and language options: "English", "Русский", "O'zbek", and "Yangi versiya". The search input field contains the word "kutubxona". Below the input field, there are buttons for "IZLASH" and "TOZALASH". A message states: "Belgilangan shartlar bo'yicha natija mavjud emas. Balki izlaganingiz quyidagilardir". Below this message, a list of search results is shown: "kutubxona", "kutubxonachi", "kutubxonalararo", "kutubxonachilik", "kutubxonashunos", and "kutubxonashunoslik".

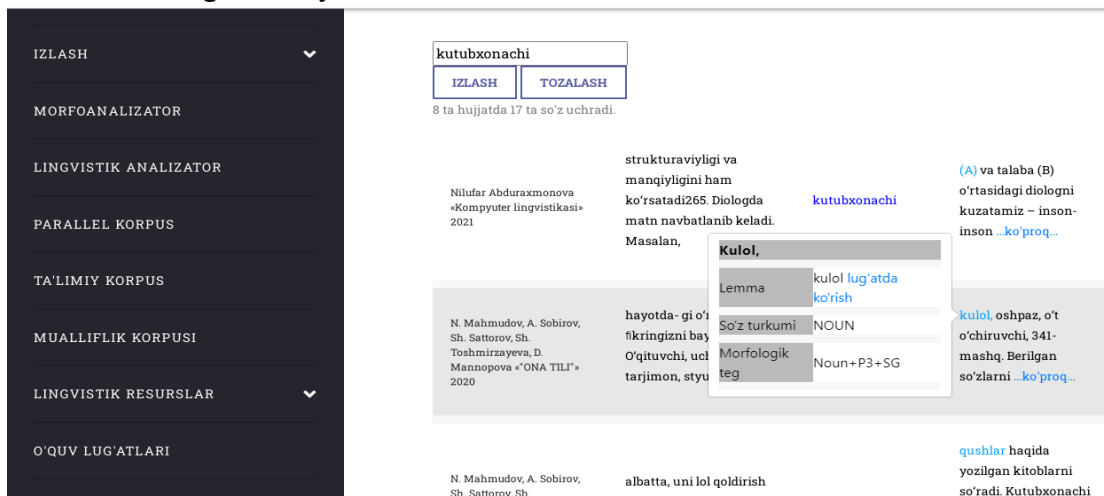
(asosan, Grammatik belgilari) bo'yicha toppish imkoniyati mavjud.

3) konkordans bo'yicha qidirish – n-gram left  $[W_1+W_2+W_3+L_t]$ / n-gram right  $[L_t + W_1 + W_2 + W_3 + W_n]$ . Bunda  $W$  – so'z,  $W_1$  qidirilayotgan so'zga birish masofadagi n-so'z,  $L_t$  – [lemma + Grammatik kategoriya]



The screenshot shows the search interface of the O'zbek Tili Korpusi with search filters "USLUB", "DAVR", "TOMON", and "GRAMM" selected. The search input field contains the word "boshi". Below the input field, there are buttons for "IZLASH" and "TOZALASH". A message states: "45 ta hujjatda 833 ta so'z uchradi.". Below this message, there are three concordance results. Each result shows the source text, the concordance context, and the concordance type. The first result shows the source text "Tora MIRZAYEV, Asqar MUSOQULOV, Bahodir SARIMSOQOV «O'ZBEK XALQ MAQOLLARI» 2005", the concordance context "bor bo'lsa, er xor Er bor bo'lsa, el xor El", and the concordance type "boshiga". The second result shows the source text "Tora MIRZAYEV, Asqar MUSOQULOV, Bahodir SARIMSOQOV «O'ZBEK XALQ MAQOLLARI» 2005", the concordance context "bo'lsa, el xor El boshiga kun El yaratgan er El", and the concordance type "boshiga". The third result shows the source text "Tora MIRZAYEV, Asqar MUSOQULOV, Bahodir SARIMSOQOV «O'ZBEK XALQ MAQOLLARI» 2005", the concordance context "El boshiga kun El yaratgan er El boshiga tushgani Er", and the concordance type "boshiga".

Konkordans qidirishda qidiruv tizimidagi soʻzning oʻzidan oldingi yoki keying keluvchi birliklar bilan birikish holati n-gram boʻyicha koʻrsatiladi.

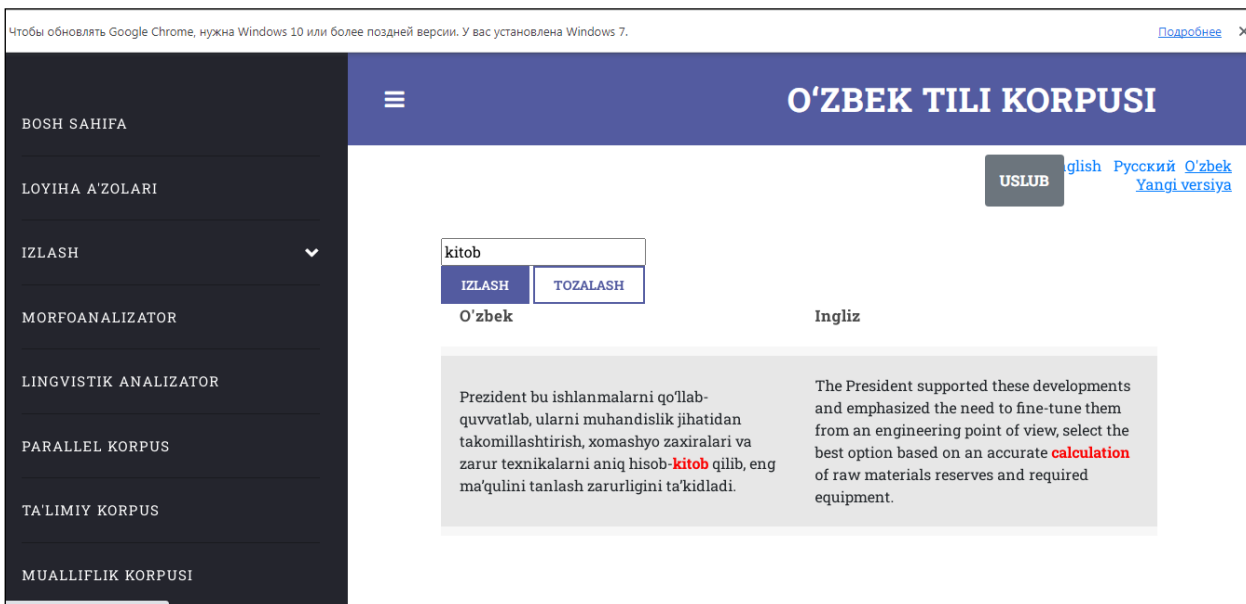


The screenshot shows a search interface for the word "kutubxonachi". It includes a search bar with the word entered, buttons for "IZLASH" (Search) and "TOZALASH" (Clear), and a dropdown menu with options like "IZLASH", "MORFOANALIZATOR", "LINGVISTIK ANALIZATOR", "PARALLEL KORPUS", "TA'LIMIY KORPUS", "MUALLIFLIK KORPUSI", "LINGVISTIK RESURSLAR", and "O'QUV LUG'ATLARI". The search results display the word "kutubxonachi" with its frequency (8 ta hujjatda 17 ta so'z uchradi) and a list of occurrences with their contexts. A tooltip for the word "Kulol" shows its lemma, word type (NOUN), and morphological tag (Noun+P3+SG).

Matndagi ixtiyoriy soʻzning lemmasi, uning turkumi va FST metodidan olingan morfologik annotatsiyasi oynada oʻz ifodasini topadi

4) soʻz birikmasi boʻyicha qidirish – [W + W + W<sub>n</sub> + Grammatik kategoriya]. Ushbu qidiruvda turgʻun birikmalarni ajratishda qoʻl keladi. Masalan, oʻrmon qiroli, dala malikasi kabi perifrazalarda gab boʻlagi vazifasida grammatik qoʻshimchalar bilan kelish imkoniyati hisobga olinadi.

Ushbu tizim ham lotin, ham kirill grafikasi boʻyicha qidirish imkoniyatiga ega. Oʻzbek tili elektron korpusining ichki korpusi parallel korpus boʻlib, unda oʻzbekcha-inglizcha va inglizcha-oʻzbekcha parallel matnlarning tarjimon xotirada segmentlangan birliklar boʻyicha toppish imkoniyati mavjud.



The screenshot shows the "O'ZBEK TILI KORPUSI" interface. It features a search bar with the word "kitob" entered, buttons for "IZLASH" and "TOZALASH", and a dropdown menu with options like "BOSH SAHIFA", "LOYIHA A'ZOLARI", "IZLASH", "MORFOANALIZATOR", "LINGVISTIK ANALIZATOR", "PARALLEL KORPUS", "TA'LIMIY KORPUS", and "MUALLIFLIK KORPUSI". The search results display the word "kitob" with its frequency (8 ta hujjatda 17 ta so'z uchradi) and a list of occurrences with their contexts. A tooltip for the word "Kulol" shows its lemma, word type (NOUN), and morphological tag (Noun+P3+SG).

Korpus qidiruv tizimi sifatida soʻzning bexato ekanligini tasdiqlash yoki uning xatoligini unga yaqin soʻzlar bazasidan qidirish korpus texnologiyasi uchun muhim sanaladi [Abdurakhmonova, N. 2019, 2021, 2022].

Korpus lingvistikasining rivojlanish evolyutsiyasi shuni koʻrsatadiki, ushbu fan kompyuter lingvistikasi bilan uzviy bogʻliq. Chunki u yoki bunda yaratilgan ilmiy yutuqlar bir-birini

to'ldirishga, imkoniyatlardan foydalanishning optimal usullaridan foydalanishga omil bo'lib xizmat qiladi. Korpuslar uchun berilgan ta'riflar ichida Webster lug'atida tilga oid ta'rifida "bilimlar yoki dalillar to'plami; tilni tavsiflovchi tahlil uchun foydalanilgan so'zlar to'plami" izohi berilgan, boshqa manbalarda korpus (ko'pligi - corpora) yozma matnlar yoki yozib olingan og'zaki nutqning transkripsiyasidan tuzilgan lingvistik ma'lumotlar to'plami bo'lib, asosiy maqsadi tilda mavjud farazlarni tasdiqlash deya, ta'riflanadi. Korpus haqidagi ta'riflar va farazlar biroz farq qilsa-da, ammo, barcha fikrlar birlashib, korpus til birliklarining xususiyatlarini aniqlash maqsadida qidiruv dasturiga bo'ysundirilgan matnlar majmui, tabiiy tildagi elektron shaklda saqlanadigan yozma yoki og'zaki, kompyuterlashtirilgan qidiruv tizimiga dasturiy ta'minot asosida joylashtirilgan matnlar jamlanmasi ta'rifini shakllantiradi.

### Foydalanilgan adabiyotlar ro'yxati:

1. Мирзиёев Ш. Миллий тараққиёт йўлимизни қатъият билан давом эттириб, янги босқичга кўтарамиз. – Тошкент: Ўзбекистон, 2017. – Б.168.
2. Abdurahmonova N. O'zbek tili elektron korpusining kompyuter modellari. Monografiya. – Toshkent, 2021.
3. Abdurakhmonova, N., Tuliyev, U., Ismailov, A., & Abdurahobov, G. (2022). Uzbek electronic corpus as a tool for linguistic analysis. In Компьютерная обработка тюркских языков. TURKLANG 2022 (pp. 231-240).
4. Abdurakhmonova, N. Z. Q., & Urazaliyeva, M. Y. (2022). O'zbek tili elektron korpusida (<http://uzbekcorpus.uz/>) og'zaki matnlar korpusini yaratishning nazariy va amaliy masalalari. Academic research in educational sciences, 3(3), 644-650.
5. Mengliev, D., Barakhnin, V., & Abdurakhmonova, N. (2021). Development of intellectual web system for morph analyzing of uzbek words. Applied Sciences, 11(19), 9117.
6. Abdurakhmonova, N. (2019). Dependency parsing based on Uzbek Corpus. In of the International Conference on Language Technologies for All (LT4All).
7. Agostini, A., Usmanov, T., Khamdamov, U., Abdurakhmonova, N., & Mamasaidov, M. (2021, January). Uzwordnet: A lexical-semantic database for the uzbek language. In Proceedings of the 11th Global Wordnet conference (pp. 8-19).
8. <https://www.merriam-webster.com/dictionary/corpus>
9. Crystal D. An Encyclopedic Dictionary of Language and Languages. - Oxford, 1992.