



# ESTABLISHMENT OF A NATIONAL CORPUS THE UZBEK LANGUAGE IS A REQUIREMENT OF A NEW ERA



**Guli Toirova Ibragimovna**  
Doctor of Philosophy, Associate Professor  
Bukhara State University

2021

# **Keywords:**

corpus, coding, transformation and graphic analysis, technological process, automatic marking, metadata, text information of a large array, morphological marking, spelling module, morphological module, linguistic module, phrase modules, word algorithm, formula algorithm, tabular algorithm, graphic algorithm.

In world linguistics, by the second decade of the twenty-first century, the creation of language corpora on the Internet is the main means of maintaining a particular language, expanding its field of research and demonstrating language skills. In particular, computer technology, which is the great invention of the twentieth century, opens the door to a wide range of possibilities for linguistics as well as other fields and poses important challenges for the computer language, the emergence of computational linguistics is critical to the success of natural languages. In global language studies, the study of linguistic language modeling, the development of algorithms for lemming words and tags, as well as the electronic use of oral and written monuments, samples of spiritual heritage created in a specific language, in order to increase the use of national and cultural heritage. Particular attention is paid to the processing of information using computer technology, the development of the necessary methodological and software for the introduction of information resources, the development of the language corpus on the Internet and, based on this, the scientific and theoretical aspects of the national language.

On the technological process of building a body that provides the stages of the technological process VV

Rykov, Yu.N. Marchuk, I. Melchuk, Sh. Khamroeva

1. Лексик таҳлил – матнни гаплар ва сўзларга ажратиш.
2. Сўзларнинг морфологик таҳлили – сўзлар контекстини (маъно устуворлигини) инобатга олган ҳолда нутқ қисмини, келишигини, турини (родини) ва бошқа грамматик белгиларини аниқлаш.
3. Лемматизация – сўзни дастлабки шаклига (леммега) келтириш.
4. Синтактик таҳлил (dependency parsing) – гаплардаги сўзлар боғланишини аниқлаш, эга ҳамда кесимни излаш ва ҳ.к.з.
5. Соддалаштирилган синтактик таҳлил (chunking) – гапларнинг эга, тўлдирувчи ва ҳол бўйича гурухларга ажратилиши.

# Компьютер технологиясида ишлатиладиган интерфейс қуидаги турларда бўлади:

*Визуал.* Мониторда намойиш этиладиган визуал тасвирлар ёрдамида маълумотларни узатадиган стандарт компьютер интерфейси.

*Имо-ишора.* Қоида тариқасида, у телефонлар ёки планшетлар учун интерфейс бўлиб хизмат қиласи. Кўпгина ҳолларда, бу тизимни бошқарадиган одамнинг бармоқларининг харакатларига жавоб берадиган ва ҳар бир аниқ ҳаракатга маълум даражада жавоб берадиган сенсорли панел. Оддий визуал интерфейснинг соддалаштирилган версияси деб аташ мумкин.

*Овоз.* Ушбу турдаги интерфейс нисбатан яқинда пайдо бўлди. Овозли буйруқлар ёрдамида тизимни бошқариш имкониятини беради. Тизим, ўз навбатида, фойдаланувчи билан мулоқот орқали ҳам жавоб беради. Энг қизиғи шундаки, замонавий технологиялар бизга нафақат телефонлар ёки компьютерларнинг овозини, балки майший техника ва ҳатто бортли компьютерларнинг овозини бошқаришга имкон беради.

This is the search window of the case, i.e. the interface.  
"Body search" - this is where you can search for a word or phrase.

The screenshot shows a web-based search interface for the 'Uzbek Tilining Milliy Korpusi'. The top navigation bar includes tabs for 'Корпусдан қидириш' and 'About.aspx'. The main menu on the left lists various search functions: Асосий, Корпусдан қидириш, Токенайзер, Лемматайзер, Разметкалаш, Частотали лугат, Корпус нима?, Корпус яратиш учун асослар, Семантик разметканинг теглаш моделлари, Морфологик разметкалашнинг лингвистик модел ва теглари, Синтактик разметканинг теглаш моделлари, Корпус статистикиаси, Таҳлиллар, Онлайн лугат, and Лойиҳа ҳақида. The central search area is titled 'Корпусдан қидириш' and 'Сўз ёки матн қидириш'. It features a search input field with buttons 'Қидирув' and 'Тозалаш'. Below this is a section titled 'Лексик-грамматик қидириув' containing fields for 'Сўз' (with input 'эшит'), 'Грамматик ҳолати' (empty), 'Омоним' (empty), 'Кўлланилиши' (empty), 'Пароним' (empty), 'Вариантни' (empty), 'Синоними' (empty), 'Антоними' (empty), 'Услуби' (empty), and 'Кўлланилиш даври' (empty). Another search input field with 'Қидирув' and 'Тозалаш' buttons is located below. At the bottom, there is a section titled 'Синтактик қидириув' with fields for 'Содда гап' (empty), 'Кўшма гап' (empty), 'Уюшган гап' (empty), 'Даррак гап' (empty), and 'Буйруқ гап' (empty). A Microsoft Word document titled 'Документ1 - Microsoft Word (Сбой активации продукта)' is open in the background. The system tray at the bottom right shows icons for RU, battery level, and date/time (12:01, 07.05.2020).

# In this case, the result will look like this :

Корпусдан қидириш - My ASP.NET

localhost:6198/About.aspx

## ЎЗБЕК ТИЛИНИНГ МИЛЛИЙ КОРПУСИ

### Асосий

#### Корпусдан қидириш

- Токенайзер
- Лемматайзер
- Разметкалаш
- Частотали лугат
- Корпус нима?
- Корпус яратиш учун асослар
- Семантик разметканинг теглаш моделлари
- Морфологик разметкалашнинг лингвистик модел ва теглари
- Синтактик разметканинг теглаш моделлари
- Корпус статистикаси
- Таҳлиллар
- Онлайн лугат
- Лойиҳа ҳақида

#### Сўз ёки матн қидириш

маъно

Кидириш Тозалаш

Сизнинг сўровининг бўйича 0 та сўз топилди!

Сўз	Маъноси
мақол	Хаётый тажриба асосида халқтомонидан яратилган, одатда панд-насиҳат мазмунига эга бўлган ихчам, образли, тугал маъноли ва ҳикматли ибора, ган.
сўз	1 тлш. Ўз товуш қобигига эга бўл-ган, нарса-ходисалар, жараёнлар, шахслар, белги ва миқдорларни, хусусиятларни, харакат ва ҳолатни, алоқа ва муносабат-ларни номлаш учун хизмат қиласиган, мус-тақиул луғавий маънога эга бўлган, шу-нингдек, турли грамматик маъно ва вазифаларда кўпландиган энг мухим тип бир-лиги; гапнинг курилиш материали.
учун	Сабаб, мақсад, атапланник каби маъноларни билдиради, шундай маъноли муносабатларни кўрсатади.
факат	1 рвш. Биргина, ягона, ёғиз.2 зидл. боғл. Чеклаш маъносини ифода-лайди.
экан	Феъл, шунингдек, бошка туркумга оид сўзларга бириниб, ке-йин билганилик, шарт ва б. модал маъно-парни билдиради.
энди	1. Маънони кучайтиради, таъ-қидлади. 2. Ҳаракат-ходисаларнинг давомлилиги, кетма-кетлигига сўнгтисини, навбатда-гисини билдиради. 3.) Ҳаракат-ходи-саннинг ҳозиринга юз берини, бошланишини, юз берганига (бўлганига) ҳали ҳеч қанча вақт бўлмаганини билдиради.
бир	1 1 раками ва шу рақам билан ифодаланган энг кичик сон, миқдор. 2 Ҳисоб, миқдор, ўлчам, чама маъно-ларини англатувчи сўзлар билан келиб, предметларнинг миқдори, ўтчамини анг-латади. 3 Баъзи отppardан олдин келиб, но-аниклик тушучасини ифодалайди ва «қан-дайдир», «аллақандайдир», «аллақайдай» каби маъноларни билдириб келади.
бирига	1 Бир бўлиб, қўшилиб, бир-галиқда, ҳамкорлика. 2 Бир вақтда, баравар. Углим эрталабмен башан бирга туради. 3 Билан кўмакчисидан кейин келиб. «қаторида, бир қаторда, шунингдек» маъно-ларини англатади.
бўлиб	1 Бўлмоқ 2 От туркумидаги сўзлар билан бириниб, «вазифасида», «сифатида», «ўхшаб» маъно-сида ишлатилади
кол	маъно - Ҳаётый тажриба асосида халқтомонидан яратилган, одатда панд-насиҳат мазмунига эга бўлган ихчам, образли, тугал маъноли ва ҳикматли ибора, ган.

#### Лексик-грамматик қидирив

Сўз ? Документ1 - Microsoft Word (Сбой активации продукта) ТИ ? Омоним ?

Нет подключения - Есть доступные подключения

RU

11:50  
07.05.2020

Корпуснинг “Лексик-грамматик қидириув” тутмасидан сўзнинг қуидаги хусусиятларини билиб олиш мумкин: Сўзнинг маъноси, грамматик ҳолати, синонимлиги, омонимлиги, паронимлиги, антонимлиги, сўзнинг варианти, сўзнинг қўлланилиш даври, сўзнинг қўлланилиш услуби кабилар.

The screenshot shows a web-based application for searching the Uzbek Corpus. The main header reads "ЎЗБЕК ТИЛИНИНГ МИЛЛАЙ КОРПУСИ". On the left sidebar, there are several menu items: Асосий (Main), Корпусдан қидириш (Search from corpus), Токенайзер (Tokenizer), Лемматайзер (Lemmatizer), Разметкалаш (Markup), Частотали лугат (Frequency dictionary), Корпус нима? (What is the corpus?), Корпус яратиш учун асослар (Basics for creating a corpus), Семантик разметканинг теглаш моделлари (Semantic markup model), Морфологик разметкалашнинг лингвистик модел ва теглари (Morphological markup model and linguistic models), Синтактик разметканинг теглаш моделлари (Syntactic markup model), Корпус статистикаси (Corpus statistics), Таҳлиллар (Analysis), Онлайн лугат (Online dictionary), and links to uzbekcorpora.uz and Проигрыватель Windows Media.

The central part of the interface is titled "Лексик-грамматик қидириув" (Lexical-grammatical search). It contains two main sections: "Корпусдан қидириш" (Search from corpus) and "Синтактик қидириув" (Syntactic search).

**Корпусдан қидириш:**

- Сўз ёки матн қидириш:
  - Кидирив
  - Тозалаш
- Лексик-грамматик қидириув:
  - Сўз ?  
Кўлланилиши ?  
Синоними ?  
Кўлланилиш даври ?
  - Грамматик ҳолати ?  
Пароним ?  
Антоними ?  
Дарак гап ?
  - Омоним ?  
Варианти ?  
Услуби ?  
Сўроқ гап ?
- Кидирив
- Тозалаш

**Синтактик қидириув:**

- Содда гап ?  
Дарак гап ?
- Қўшма гап ?  
Сўроқ гап ?
- Уюшган гап ?  
Буйруқ гап ?

At the bottom, there are system status icons and a date/time stamp: RU, 12:00, 07.05.2020.

ASP.NET

ut.aspx

## Корпудан қидириш

Сүз ёки матн қидириш

Кидириу  
 Тозалаш

### Лексик-грамматик қидирув

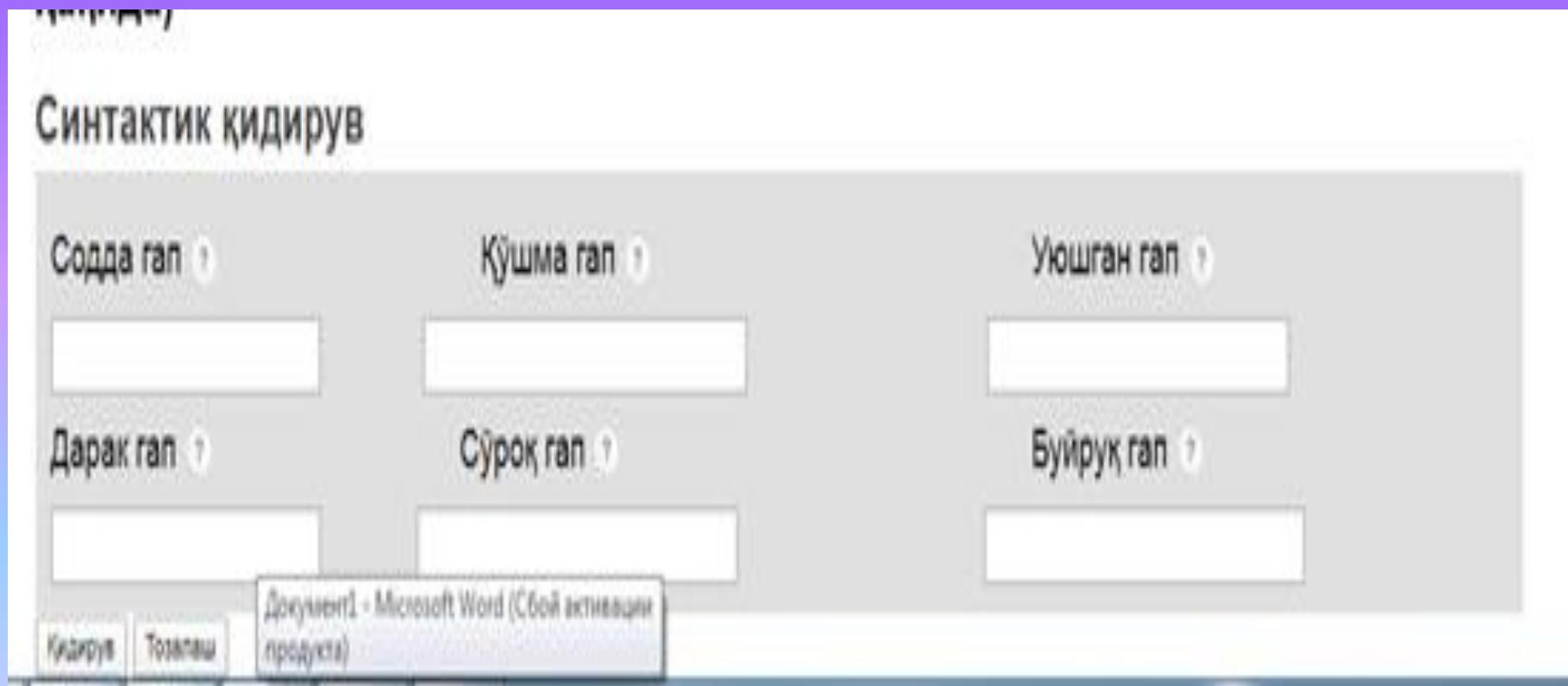
Сүз ? <input type="text" value="ингичка"/>	Грамматик ҳолати ? <input type="text" value=" [сиф.], [ш. ЛМГ]"/>	Омоним ? <input type="text" value="Йүк"/>
Құлланилиши ? <input type="text" value="чегараланмаган"/>	Пароним ? <input type="text" value="Йүк"/>	Варианти ? <input type="text" value="Йүк"/>
Синоними ? <input type="text" value="майин"/>	Антоними ? <input type="text" value="қалин"/>	Услуби ? <input type="text" value="бетараф"/>
Құлланилиш даври ? <input type="text" value="замондош"/>		

**ИНГИЧКА - 1, Кундаланг кесими меъёр-дан кичик; 2 Чийилдоқ, үткір (товуш ҳақида)**

### Синтактик қидирув

Содда гап ? <input type="text"/>	Құшма гап ? <input type="text"/>	Үюшган гап ? <input type="text"/>
-------------------------------------	-------------------------------------	--------------------------------------

# Syntactic search button



# According to the structure of the sentence, simple sentences are divided into the following groups:

Бир таркибли содда ёйиқ гаплар	Икки таркибли йиғиқ содда гаплар	Икки таркибли ёйиқ содда гаплар
-Бұлсанғ-чи! Ҳали замон қун ёйынан кетади. Далага чиқышы керак.	Инсон – Ҳазрат! Инсон – Қабир! Бу гапни аммо уққаным ішік. На падарим еа на онамдан.	Жиіннингизни олиб кетганим келдім (шахси аник).
Кече дүнөнің қырғынлы туны экан.	Умматали индамади. Үрнидан түриб у ёқдан-бу ёққа юра б	-Гулзорни бузиб ташладилар (шахси ноаник).
Адабини берши керак зди. Ең күнчилек орасида хүрлаб сұқамизми?	Хәсіп иссік. Күн бүйін юғуриб-еңіб қарчаган одамлар еа отлар ҳализамон chanқаб сүв қидиришега мүжаррар.	Бугунги ишни әртага күйма (Макол) (Шахси умумлашган).
Соатти графикдан түгри фойдаланылсın.	Айбы кантта. Үни шүңдай қолдырыб бўлмайди.	Кула-кула ўлибмиз. Буни Киғояхоннинг ўзи боллаб айтади.
Үйланишни күт көчкітірмаслик керак, иним. Ёшлик хам ғанымат (Сўзлашувдан).	Хонзода бевгим шигирма беш ёшда. Айни узатыладиган пайтлари.	Начора, шхитиёрги. Нима қылмоқчи бўлаётганингни ўзинг ҳам тушукмасанг керак.
Эй, кіхна Шарқимизнинг ҳеч сұнмас машъали.	Китоб иккі нозу јұттыз беттілек экан. Жұда қызықарлы.	Ҳамма газета, радио еа телевидение сизнинг шхитиёргиздә!
Үқтін-үқтін шаҳарға түшиб Думасига етиб бориш керак.	Бақор сенікішір.	Әзина! Ганциринг! (жозирғи замон, буйруқ-истак майли)
Ҳамманғи қатор құйиб битта-битта пешонанедан отиши керак.	Гап шу. Ашрафжон шүңдай дөди-да, үрнидан түриб кетди, сүне аңчагача түғон атрофида айланыб юрди.	... Ҳа, ҳа, хунуксан! (жозирғи-көласи замон, аниқлік майли)
- Түғри, ийк. Демак, ойшыларнинг устларидан ҳеч қаерге шикоят қылыш бўлмайди.	Буен юрсақ, әртага яхши ўлларга чиқамиз. Йўл анча.	-Қисқаси, баҳтиёрсизлар. (жозирғи-көласи замон, аниқлік майли)
	Низомжон индамади. Юрагида тўлиб-тошган гапларни	

# O'ZBEK TILI MILLIY KORPUSI



Alisher Navoiy  
1441-1501

ROMANLARI  
HIKOYALARI  
PYESA  
PORTRET  
FELYETON

- Alisher Navoiy hayoti va ijodi
- O`zbek mumtoz adabiyoti

## FOTO ALBOM

+ Buyuk siymolar

## TAQDIMOT

+ Taqdimot



## Korpus nima...

Korpus — til birliklarining xususiyatlarini aniqlash maqsadida qidiruv dasturiga bo'yusundirilgan matnlar majmui, tabiiy tildagi elektron shaklda saqlanadigan yozma yoki og'zaki, kompyuterlashtirilgan qidiruv tizimiga dasturiy ta'minot asosida joylashtirilgan on-line yoki off-line tizimda ishlidaydigan matnlar jamlanmasi[1]. U elektron shaklda mavjud bo'lgan matnlarni, hujjalarni ma'lumotlarni qayta ishlash, ularni avtomatik tahlil, ya'ni morfologik sintaktik va semantik tahlil qilish, morfologik analiz va sintez qilish bilan birga, yarim avtomatik rejimda ma'noni buzmadsan ishonchli nutq materialini moslashtirish darajasini tekshirish ushu sohadagi til bilan qiziqishni oshirishga xizmat qiladi. Korpus — tilini tadqiq etish, til o'rganish, lug'at tuzishda eng zamonaevi, keng imkoniyatlari dasturlashtirilgan tizimdir. Korpus - millionlab so'z kontekstlariga havola qiluvchi, maxsus qidiruv tizimi asosida ishlidaydigan elektron matnlar yig'indisi.U internet tizimidagi elektron kutubxonasi, lug'at va lingvistik grammatikadir. Korpus - tabiiy (real) tilning elektron shaklda bo'lib, uni to'ldirish, tuzatish, tahrir qilish imkoniyatiaga ega[2]. Korpus bu so'zlar, so'z birikmaları, grammatic shakllari ma'nosini ma'lum bir izlash tizimi orqali topishning elektron shakldagi matnlar to'plamidir. Korpuslarning har xil turlari mavjud. Masalan, bir muallif korpusi[3], bir kitob korpusi (jumladan dastlabki korpuslar "Bibliya" uchun qilingandir). Ma'lum bir til Milliy korpusi shu tilning xayotining barcha qirralarini, janrlarini, usullarini, xuduydi va ijtimoiy variantlarini o'zida ifoda etadi. Masalan, o'zbek til Milliy korpusini Internetda amalga qo'yilsa, 140 min so'z ishlatlardan iborat shu tildagi barcha turdag'i matnlarni o'zida mujassamlashtira oladi. Kelajakda o'zbek til Milliy korpusi o'zida 300 mln so'z ishlatlarni jamlashi lozim bo'ladi. Boshqa milliy tillar korpuslari kabi o'zbek til Milliy korpusi xam ikkita muxim xususiyatga egadir. Birinchidan u ancha salobatlari va barcha yo'nalishlar uchun (turli badiy janrlar: publisistik,o'quv,ilmly, ish yuritish, nutk - so'zlashuv, shevalar va boshqalar uchun) bir xilda muvofiglashgandir. Bu matnlar o'z davrida tildagi ta'sir kuchiga qarab proporsional ravishda korpusda jamlangan bo'ladi. Ikkinchidan korpus matnlarga xos bulgan alohida jihatlarini qo'shimcha ma'lumot sifatida o'zida aks ettiradi (masalan ma'lum bir jixatdan belgilanishlar va annotatsiyalar). So'zlar va boshqa birikmalar belgilanishi - korpusning bosh xarakteristikasıdir (tasnifidir); U korpusni hozirda Internetni to'ldirib yuborgan boshqa matnlar kollektsiyalari va kutubxona matnlardan farglaydi. Matnlar belgilanishi qanchali boy va turli tuman bo'lsa korpusning ilmiy ahamiyati shuncha yuqori bo'ladi. O'zbek til Milliy korpusini belgilanishlarning uch turi e'tiborga olinish kerak bo'ladi: metamatnli (u matnni muallif qarashi, janri va boshqa xususiyatlar bo'yicha to'liq ifodalaydi), morfologik va semantik belgilanishlar (morphologik va semantik belgilanishlar matnni emas balki alohida bir so'zni tasniflab beradi). ADABIYOTLAR: 1. Grudeva V.A. Korpusnaya lingvistika. Uchebnoe posobie. 2-e izd., stereotip.- M.: Flinta, 2012. – 165 s. 2. Zaxarov V.P. Korpusnaya lingvistika: Uchebno-metod. posobie. – SPb., 2005. – 48 s. 3. Zaxarov V.P., Bogdanova S.Yu. Korpusnaya lingvistika.– Irkutsk: IGLU, 2011. – 161 c. 4. Tolrova G.I. Importance of Interface in Creating Corpus. International Journal of Recent

## SO'Z QIDIRISH

Hikoya nomi:
Kartina
<b>Qidiruv</b>

## LEKSIK-GRAMMATIK QIDIRUV

- So'z
- Grammatik holati i
- Omonim
- Qo'llanilishi
- Paronim
- Varianti
- Sinonimi
- Antonimi
- Uslubi
- Qo'llanilish davri

**Qidiruv**

## SINTAKTIK QIDIRUV

- Gapning tuzilishiga ko'ra turlari
- Gapning maqsadiga ko'ra turlari

**Qidiruv**

## MUALLIFLAR

**Mualliflar**

## DASTUR HAQIDA

Ushbu korpus qidiruv natijasi shu hujoya matni asosida amalga oshiriladi. Korpus



# Conclusion

1. Corpus linguistics is the most advanced branch of linguistics, and corpus is a necessary tool for linguists; oral, written monuments are a source of information reflecting the national and cultural heritage. A corpus is a set of texts to be searched for, and a well-defined corpus serves as a stable linguistic base to ensure the effectiveness of linguistic research. As an artificial intelligence product, the linguistic corpus includes an electronic dictionary, translation portal, terminology database, virtual (electronic) library, electronic government, electronic publications, electronic textbooks and manuals. Linguistic electronic sources, which are a product of artificial intelligence, are considered raw materials for creating a certain linguistic corpus.

2. The creation of the National Corpus is carried out in two stages: determination of the list of sources and digitization of texts (transformation into computer form). Its technological process consists of: creating a dictionary of repetitions of lexemes and word forms based on selected texts; view the text for any unit of the received dictionary of repetitions; divide a graphic word into syllables and compose a dictionary of repetitions of syllables; sorting word resources; simultaneous processing of an unlimited number of files; create text corpora with external symbols; calculation of statistical data for the corpus of created texts and individual texts included in the corpus; work with source texts in txt, doc i rtf format, automatic encoding setting, etc.

3. The most effective standards for encoding corpus data have been selected. The presentation of its data is based on an SGML / XML text layout. When encoding, lexical information is adapted to HTML / XML rules. The texts selected for the National Corpus are taken from various sources and presented in different formats: plain text, HTML, RTF, PDF.



**Thank you for your  
attention !!!**