

## Татар теленең язма корпусы

21 нче гасыр башында Интернет битләрендә, телчеләр телендә “корпус”, “милли корпус” кебек терминнар күзәтелә һәм ишетелә башладылар. 1990 елда инглиз теленең милли корпусы, соңнанрак нимес, француз, итальян телләренең милли корпуслары төзелеп дөньяга таралды. 2007 елларда рус теленең милли корпусы төзелеп бетүе мәгълүм булды. Шул вакытлардан бирле төрле илләрдә милли корпуслар туа торды. Нәтижәдә компьютер һәм корпуслар ярдәмендә тел өйрәнүче фән тармагы “Корпуслы лингвистика” (Корпусная лингвистика) барлыкка килде, телчеләр арасында 21- нче гасырны милли корпуслар гасыры дип тану гадәتكә керде.

Төрки телләрдә күрсәтелгән янарышлар соңгы елларда гына күзгә чалына башладылар. Интернет хәбәр итүенчә, төрек теленең ике миллион сүздән торган тамгаланган (разметка ясалган) корпусы төзелгән. Кызганычка каршы, ул корпус белән танышу өчен регистрация үтү сорала (коммерческий корпус булырга охшый). Казах, кырым татарлары телләре өчен корпуслар булуы турында әйтелә, әмма аларның күләме, сыйфатлары хакында мәгълүматлар юк.

### Нәрсә ул корпус?

Гомумиләштереп әйткәндә, корпус - ул күп сандагы текстлар жыелмасы нигезендә төзелгән һәм электрон формада саклана торган сүзлек–белешмә. Корпуста, ике телле сүзлекләргә хас булган, сүzlәрнең тәржемәсе дә аңлатмалы сүзлекләргә хас булган төшендерүләр дә китерелми. Аның вазыйфасы телне электрон формада мөмкин кадәр тулы һәм барлык үзенчәлекләре белән гәүдәләндерү.

Моңа чаклы телләр ике формада күзәтелеп килде. Бүгенге компьютерлар телне тагы бер формада – электроннар жирлегендә яшәү мөмкинлеге белән тәэмин ителәр. Мәгълүм ки, бүгенге компьютерларның хәтерә телне язма текстлар буларак та, сөйләм текстлар буларак та саклауга; белемә исә (программалар системасы) аларны бер формадан икенче формага күчәргә житәрлек. Гадиләштереп әйткәндә, алар шәхес куйган сорауга сөйләм телендә дә, язма телдә дә җавап бирергә сәләтлеләр.

Шәхес куйган сораулар нинди булмаска мөмкин? Нәкъ менә шуның өчен дә тел компьютер хәтерендә тулы һәм барлык үзенчәлекләре белән гәүдәләнергә тиеш була.

Бер уңайдан тел төшенчәсен дә ачыклап үтик. Бөекбритания энциклопедиясендә: “Тел – ул аерым халык, милләт тарафыннан фикер гәүдәләндерү, тапшыру максаты белән файдаланучы барлык сүз һәм сүзтезмәләрнең жыелмасы”- диелгән.

Электрон корпус шушындый жыелманы тәшкит итә дә инде. Милли корпус дәрәжәсендәге корпусларларның күләме 100 миллион сүздән (сүзкулланыштан) ким булмый. (Инглиз теленең милли корпусы 400 млн. сүздән рус теленең милли корпусы 300 млн. сүздән тора). Милли корпуслар

язма тел, сөйләм тел, жирле һәм социаль диалектларны чагылдырырга тиеш булалар.

Корпусларга хас тагы бер сыйфат – ул, нинди төр лингвистик мәсьәләләрнең корпус ярдәмендә чишелү мөмкинлеген күздә тотып төзелә. Шул сәбәпле корпус, лингвистик мәгълүматлар белән генә түгел, күздә тотылган мәсьәләләрне чишү программалары белән дә тәмин ителә.

Кайсы төрдәге лингвистик мәсьәләләрне чишүгә юнәлтелгән булуларына карап корпуслар күп кенә типларга бүленәләр:

1. Милли һәм милли булмаган корпуслар. Милли корпуслар күп төрле мәсьәләләр чишүгә мөмкинлек тудыралар һәм, алда әйтелгәнчә, күләм ягыннан 100 миллион сүздән ким булмыйлар. Милли корпусларга хас икенче бер сыйфат - ул репрезентатив һәм сбалансировланган булырга тиеш.

2. Нинди текстлар жыелмасын тәшкит итүләренә карап корпуслар язма (письменный), сөйләм (устный) һәм катнаш (смешанный) булалар.

3. Телләр типологиясе, тәржемә (машина тәржемәсе) кебек мәсьәләләргә мөнәсәбәтле корпуслар бертелле (однойзычные), икетелле (двухязычные) һәм күптелле (многоязычные) булырга мөмкиннәр һ.б.

Әйтергә кирәк, корпусларга хас төрлелек әйтелгәннәр белән генә чикләнми. Аерым фәннәр телен, гасырлар дәвамында конкрет бер телнең үзгәреш динамикасын чагылдырган төрле чор корпусларының һәм хәтта ябык корпусларның булуы да мәгълүм.

Без төзегән язма корпуска нинди сыйфатлар хас?

Шушы елның апрель башында Татар теленең язма корпусы төзелеп Интернетка куелды. Исеменнән күренгәнчә, сөйләм теле корпуста юк. Шуның өчен аны милли корпус дип әйтеп булмый. Ә менә язма телебезне ул житәрлек дәрәжәдә тулы чагылдыра дип әйтә алабыз. Язма корпусның күләме 45 миллионнан артыграк сүзгә тигез. Корпуска кергән сүзләр эчендә төрле сүзформаларның саны 400 мең чамасында. Корпус, татар әдәби телендә киң чагылыш тапкан, публицистика, матур әдәбият һәм гуманитар фән стилиндәге текстлар жыелмасыннан тора. Корпуста публицистикага караган текстлар якинча 60%, матур әдәбиятка караган текстлар - 35%, фәнни текстлар -5% тигез.

Корпусның программалар системасы (корпус менеджеры) тел мәсьәләләре белән кызыксынган шәхескә күп төрле мәгълүматлар бирергә мөмкин. Мисал өчен:

1. Шәхес кызыксынган сүзгә табарга һәм аның татар телендә нинди ешлык белән кулланылуын күрсәтергә,

2. Табылган сүзгә алдыннан һәм артыннан килүе мөмкин булган сүзләргә (сул һәм уң контекстларны) ачыкларга һәм аларның кулланылу ешлыкларын бәян итәргә,

3. Табылган сүзгә төрле жөмлөләрдә нинди сүзләр белән семантик мөнәсәбәткә керү мөмкинлеген (семантик контекст) күрсәтергә,

4. Табылган сүзгә телдә ничек файдаланыуын дәлилләп 50 мисал – жөмлө белән дәлилләргә.

Хәзерге вакытта корпусның күләмен һәм функциональ мөмкинлекләрен киңәйтү буенча эшләр башкарыла. 2-3 айдан корпусның күләме дә, мөмкинлекләре дә артыр дип өметләнәбез.

Корпуслар төзү, аерым алганда татар теленең язма корпусын төзү нигә кирәк?

Соңгы елларда тел белемендә корпуслардан файдаланып телне өйрәнүче яңа фән - корпуслы лингвистика (Корпусная лингвистика) барлыкка килде. Корпуслы лингвистика ул компьютер һәм корпуслар симбиозы гына түгел, ул лингвистика белән социолингвистиканы (тел төзелеше белән тел кулланышын) бергә ялгап өйрәнүче фән. Корпуста чагылыш тапкан нәтижеләр татар тел структурасыннан тыш телнең файдлануын – аерым буынның сөйләмен (тел дөнъясын) дә тасвирлыйлар. Гасыр буе әдипләр телендә табылган аерым мисалларга нигезләнеп, я булмаса гәзит битләрендә күзәтелгән терминнар таянып: “Телебез байый, үсә”- дип килдек. Баксаң, телебез югалу хәлендә булып чыкты. Кулланыш сфера һәм средасы тарайган, фәнни терминнарын югалткан аралашу системасы булуы ачыкланды.

Корпусның мөмкинлекләрен күрсәтү йөзәннән берничә мисал карап үтик.

1. *Корпус* сүзе (лингвистик термин буларак) татар телендә табыламы?

Корпустагы мәгълүматлар бу сүзнен 45 миллионнан артык сүзләр жыелмасында биш мәртәбә очыравын күрсәтәләр. Әмма бу очырауларда *корпус* сүзе техник һәм хәрби терминнар буларак файдаланыла. (Сул һәм уң контекстлар составына игътибар итегез). Димәк бүгенге буын теленә *корпус* сүзе “электрон корпус” мәгънәсендә, ягъни лингвистик термин буларак күзәтелми.

2. Бүгенге буын язма татар телендә *аманәт* сүзе 158 мәртәбә, *акция* сүзе 2474 мәртәбә очыраган булып чыкты.

3. Корпуста чагылыш тапкан мәгълүматлар көтелмәгән дә булырга мөмкин. Безнең өчен түбәндәге зурлыклар көтелмәгән булып чыктылар. Корпус нәтижеләре буенча *харам* сүзе буын телендә 140 мәртәбә, *ә халәл* сүзе 3432 мәртәбә файдаланылган. Безгә калса, *харам* сүз *халәл* сүзгә караганда күберәк күзәтелергә тиеш иде. Чөнки дин, беренче чиратта, харам ризыклардан саклануга юнәлтелән тәгълимат.

4. Кызганычка каршы, корпуста 20 – 21 гасырлар чигендә дөнъя күргән техник терминнары булдыруга, булганнарын активлаштыруга караган сүзлекләрнең телгә тәэсире сизелерлек түгел.

5. Корпуста табылган мәгълүматлар татар теленә хас булган кайбер терминнарның бүгенге язмаларда русча вариантта очыравын дәлилликләр. Мәсәлән, *ваба* (*waба*) сүзе 29, *ә* аның тәржемәсен тәшкит иткән *холера* сүзе 52 мәртәбә файдаланылган. Шулай ук *бүсер* сүзе 42 мәртәбә, *грыжа* сүзе 5 мәртәбә кулланылган булып чыкты.

Татар теле электрон корпусы ярдәмендә игътибарлы шәхес бүгенге буын телендә теге яки бу сүзнен, сүзформаның барлыгы-юклығын, булган очракта аларның нинди ешлык белән кулланылуын бик тиз ачыклай ала.

Еш кына телчеләрнең үзләре ижат иткән кагыйдәләргә, ясаган искәртмәләренә әдәби әсәләрдән мисаллар эзләп утыруларын күрергә туры килде. Һәр телченең үз темасы буенча мисаллар туплавы да, картотекалар төзүе дә күпләребезгә мәгълүм. Электрон корпус аларны мондый эшләрдән азат итәчәк.

Шундый ук ярдәмне электрон корпус татар телен ана тел яки чит тел буларак укытучыларга да күрсәтергә мөмкин. Корпустан файдаланып секунд аралыгында укытучылар үзләренә кирәкле мисалларны таба, укучылар исә үзләре күрәсе килгән сүзләргә сүзформаларны күрә, аларның кулланылу формалары, язылышы белән таныша алалар.

Корпус төзелү якин арада татар теленең ешлыклар сүзлеген, аның морфологик төзелешен өйрәнү өчен мөһим булган кирәк сүзлек, татар теленә хас стильләр үзенчәлеген, аерым әдипләр ижатының стиль үзенчәлекләрен чагылдырган сүзлекләр эшләнүгә дә мөмкинлекләр тудыра.

Телчеләргә электрон формада сакланучы тел байлыгы һәм аны тиз арада күзәтеп чыгу юлы тәкъдим ителә. Русларның милли корпусы рус телен корпуслы лингвистика фәненә индергән кебек, татар теленең милли корпусы да якин киләчәктә татар тел белемен корпуслы лингвистика кысаларында өйрәнүгә китерер дип ышанасы килә.

Әйтергә кирәк, корпус техник ярдәмче генә түгел ул – татар тел дөнъясы. Бүгенге көндә татар тел дөнъясы нинди? Нигә бүген тел ияләре аңындагы төшенчәләр, тел төшенчәләр буларак, шәхес аңында гәүдәләнеш тапмый? Шәхесләр уй-фикерләрен туган телләрендә әйтә алмый газап чигәләр.

Соңгы елларда Машина тәржемәсе, Автоматик рәвештә әйтелгәнне танып торган системаларның бары тик корпуслар ярдәмендә генә чынга ашу мөмкинлегенә тәмам ачыкланды.

Әйттик, машина тәржемәсендә система *үнкә* сүзенә юлыкты ди. Бу сүз - омоним. Омнимнар, кагыйдә буларак, контекстан чыгып чишеләләр. Әмма аерым очракларда чишелеш сул һәм уң контекстлар хәленнән килмәскә дә мөмкин. Андый очракларда, гадәттә, текстның тематик үзенчәлегенә ихтибар ителә. Тематик үзенчәлек *үнкә* сүзнең семантик контекстында чагылыш тапмый калмый. Кыскасы, әгәр тәржемә ителүче *үнкә* сүзенең семантик контексты, корпустагы *үнкә* сүзенең семантик контекстына охшаш икән, ул чагында *үнкә* “легкие” мәгънәсендә тәржемә ителергә тиеш була.

Татар теленең төрле чорга караган корпусларын булдыру, аларны үзара чагыштырып телебез кичергән хәлләргә, аның үткән юлын, киләчәккә төсмерләү бүгенге чорда хыялга сыймаган эшләр түгел. Секунд эчендә кирәкле сүзгә, кирәкле форманы; галимнәр, укытучылар, укучылар сораган мисалларны табуга сәләтле корпус, алда күрсәтелгән мәсьәләләрнең көн тәртибенә куелышына, чишелешенә ышаныч уята.

Корпус Интернетка куелган. Аның адресы:

**<http://corpus.tatfolk.ru>**

**Өстәмә сүз.** Корпусны эшләү бүгенге көндә дәвам итә. Без аның күләмен 100 миллион сүзгә житкерәчәкбез. Бу вакыт эчендә корпус функциональ планда да баетылачак, ягъни бүген чишелми торган кайбер мәсьәләләрнең ике–өч айдан чишү мөмкинлеге туар.

Хөрмәтле галимнәр. Форсаттан файдаланып Сөзгә бер үтенечезне белдерәбез.

Без Сөздән татар телендә басылып чыккан фәнни хезмәтләрнең, эсәрләрегезнең электрон версиясен жиберүне үтенәбез. “Безнең хисапка баемакчы булалар”- дип уйламагыз. Без язма корпус төзеп бер тингә дә баемадык. Сөз дә корпустан түләүсез файдаланачаксыз. Бары тик эшегездә аңа таянуыгыз хакында әйтергә кирәк булыр. Анысы инде гадәти эхлак.

Сөз жибергән электрон версиягә килгәндә, тагы шуны белдерәбез. Әсәрләр корпуска жөмлэләргә бүленеп, ягъни текстлар буларак түгел, ә бәлки жөмлэләр берәмлегендә кертеләләр. Һәр жөмлэгә карата аның кайдан алынуы, авторы кем булуы хакында мәгълүматлар биреләчәк. Бәндэләргә Сөзнең ижат жимешләрен бөтен текст буларак күчереп алу мөмкин түгел.

*Мансур Сайхунов,  
Тәүзих Ибраһимов,  
Илнар Сәлимҗанов.*