

کارپس ژبپوهنه Corpus Linguistics

Bushra Khatoon

Abstract

Corpus Linguistics is an offshoot of Linguistics which emerged in 1970s and gained momentum in the present century. Computers have dramatically changed every aspect of our life and same is the case with linguistics. Like other fields of life interfaces between linguistics and computer science are also growing and are likely to be of high interest to linguists and this interest is accelerating more and more in future. Nowadays computer corpora are used for a number of purposes such as dictionary writing, lexical analysis, translation of different languages, study of semantic changes, neologisms, and learning other languages. This paper is an attempt to introduce corpus linguistics, its scope and types to the readers so as to make them familiar with this branch of linguistics and to motivate them to collect computer corpora for Pashto language.

ژبپوهنه د یوې ژبې سائنسي مطالعې ته وئیلې شي - د ژبې نه مراد څه خصوصي ژبه نه ده بلکې په عمومي توګه د ژبې سائنسي مطالعه او څیړنه - ژبپوهنې اوس د یوې غټې ونې شکل اختیار کړی دی او ډېرې څانګې او څرګنې ئې سپردې دي چې یوه نوې څانګه په کښې د کارپس ژبپوهنې په نوم پېژندلې شي - په دې مقاله کښې به د کارپس ژبپوهنې په اړه خبره کولې شي -

د کارپس ژبپوهنې تعریف مونږ داسې کولې شو:

کارپس ژبپوهنه د کمپیوټري کار پورا په مدد د ژبې تجزیه ده - دا تجزیه د کمپیوټر په مدد کولې شي - او مختلف قسمه خصوصي سافت وئیرز په کښې استعمالېږي - [1]

۱ کارپس خه ته وائي

د انگريزي کارپس لفظ د لاطيني ژبې لفظ Corpus نه اخذ شوی دے چې مطلب ئې دے Body يعني بدن، جسم - لاطيني لفظ Corpus انگريزي ته د فرانسيسي دغاړې په دوه شکلونو کښې راغلے دے -

اولنے Corpse - دا د زرې فرانسيسي Cors نه جوړ شوی دے -
او دويم Corps دا لفظ په اتلسمه پېړۍ کښې د جديدې فرانسيسي په لار راغلے دے -

دا اولنے لفظ Corps په انگريزي کښې په ديارلسمه پېړۍ کښې د Cors په شکل کښې راغلے - په خوارلسمه پېړۍ کښې هغه خپل اصلي حرف P په کښې دوباره اولگيدو- اول اول کښې ددې مطلب صرف body وه او د خوارلسمې پېړۍ په اخره کښې دا صرف Dead body يعني مړ وجود د پاره په مکمله توگه استعمالېدل شروع شو- او بل طرف ته اصلي لفظ Corpus په خپله هم په خوارلسمه پېړۍ کښې انگريزي ته راننوتو- [2]

۲ د کارپس اصطلاحي معني

Corpus واحد لفظ دے او ددې جمع Corpora يا Corpuses دے - په ژبپوهنه کښې مونږ د کارپس تعريف داسې کولے شو چې کارپس د يوې ژبې د داسې ټکرو مجموعه چې ته وئيلے شي چې د تحقيق د پاره په يوه منظمه طريقه چون کړې شوې او راغونډې شوې وي او د يوې ژبې د نمونې په توگه استعماليدے شي بلکې په جديده ژبپوهنه کښې به مونږ د کارپس تعريف داسې کوؤ چې د نمونه جاتي وېنا وي يا ليکلو شوؤ متنونو مجموعه چې په مشيني توگه لوستے شي او د گڼو ژبني معلوماتو د پاره ئې توضيح او تشریح کېدے شي - [3]

هسې کارپس هر يو متن ته وئيلے شي - ولې په اصطلاحي توگه هغه متن ته وئيلے شي چې په مشين باندې لوستے شي او په خصوصي توگه د متنونو د معلوم شمېر مجموعه چې په مشين لوستل کېدے شي - او په اصولي توگه يو کارپس د صحيح ژبنو خصوصياتو د مطالعه کولو د پاره جوړولے شي -

په ژبپوهنه کښې د کارپورا نه مراد هم کمپيوټري کارپورا دے اکر چې د کارپس ژبپوهنې اصطلاح په ۱۹۸۰ء کښې استعمال شوه خو د کارپس په بنياد د ژبې مطالعه د کا مسکي نه اول شروع شوې وه - Field linguist باؤز (Baos) ۱۹۴۰ء کښې هم دا طريقه استعمال کړې

وه او د structuralism روایت سپیر، نیومین، بلوم فیلد او د پائک هم دا طریقه په استعمال کښې وه هغه وخت به ژبپوهانو د کاغذونو په ټکرو د یوې ژبې ټکي لیکل او د پېزارونو په ډبو کښې به ئې جمع کول. وخت سره سره دا ټکي په کمپیوټرو کښې خوندي کېدل شروع شول او د مختلفو سافت وئیرو په مدد د ژبې هر اړخیزه څېړنې شروع شوې. [4]

په ۱۹۶۰ء کښې دغه شان د لاس لیکلې او د ټائپ لیکلې کارپورا په استعمال کښې وځو او د کمپیوټر لیکلې کارپورا رواج د ۱۹۷۰ء نه شروع شو په حقیقت کښې د جدیدې کارپورا غټه فائده په مشیني توگه د کارپورا لوستل دي. بل د کمپیوټر د ټولو نه غټه فائده دا ده چې ددې د سرچ، چون کولو ترتیبولو او formatting سپید ډېر زیات دے. ددې خرچ ډېر کم دے. او انساني محنت په کښې هم ډېر کم لگي. د وخت ډېر غټ بچت دے او نتیجه ئې ډېره صحیح او قابل اعتباره وي. د کمپیوټري کارپورا ژبني تجزیه نن سبا د کمپیوټري سافت وئیرو په وجه ډېره اسانه شوې ده. تېرو شلو کالونو کښې د ژبې داسې جدیدې او باریکې تجزیې د کمپیوټري کارپورا نه بغير ناممکنه وي کمپیوټري کارپورا په ژبه باندې څېړنو او تجزیو کښې ډېر غټ بدلون راوستے دے. په کارپس ژبپوهنه کښې د کمپیوټر د بې کچه استعمال د وجې دې ته نن سبا Computer Corpus linguistics هم وئیلے شي چې CCL ئې مخفف دے.

۳ د کارپس ژبپوهنې ساحه

د کارپس ژبپوهنې ساحه اوس ډېره پراخه شوې ده او د ژبپوهنې په گڼو څانگو کښې استعمالېدے شي او د ژبپوهانو کار ئې ډېر اسان کړې دے. د نظریاتي او تطبیقي ژبپوهنې ماهران د کارپورانه د خپل خپل مطلب کار اخیستے شي. د ټولنیزې ژبپوهنې ماهر د یوې ژبې د مقامي ویونکو د خبرواترو اډیو رېکارډنگ د طبقاتي فرق او لهجو مطالعه کولو د پاره استعمالولے شي. د نفسیاتي ژبپوهنې ماهر هم دغه رېکارډنگ د ژبې د لغزشونو د مطالعه کولو د پاره استعمالولے شي او د لغت لیکونکي ئې د مختلفو فقرو د تکرارېدو یعنی فریکوېنسي د پاره استعمالولے شي. گرامر پوهان ئې د ژبې د گرامري ساخت مطالعه کولو د پاره استعمالولے شي. دغه شان د ژبې اُستاد او د عدالتي ژبپوهنې ماهر ئې د خپل مقصد د پاره استعمالولے شي. [5]

کارپورا د دوو ژبو د ترجمو د مطالعې د پاره هم استعمالولے شي. د ۱۹۹۰ء نه پس د کارپس په بنیاد د دوو ژبو ډکشنرو ډېره زر وده اومونده. دغه شان د اسلوبیاتو په څېړنه کښې هم کارپورا ډېره بڼه ذریعه ده. د یو ثقافت ځانگړي اصطلاحات هم د کارپورا په ذریعه مطالعه کولے شي. [6]

کله چې یو کارپس جوړولې شي - نو راغونډ شوي متنونه په سافت Version کښې محفوظول، ترتیبول او ذخیره کول وي - د وېنا وي نمونو په نسبت لیکلې شوي متنونه محفوظول زیات اسان وي - د لیکلو شوو متنونو د پاره په عامه توګه اپتکل کرېکټر ریکګنیشن سافت وئیر Software Optical Character Recognition ضروري دے - دا سافت وئیر الیکټرانک متن فائیلو نه سکېن کوي - په داسې کارپورا کښې د ژوند د مختلفو شعبو نه متنونه شاملولې شي لکه د اخبارونو، رسالو، د سرکاري کاغذونو د تجارت د کاغذونو د خبرو اترونه نمونې اخستې شي - [7]

په یو کارپس کښې لاندیني خصوصیات په کار دي -

- په یو کارپس کښې د هغه ژبې عمومي او خصوصي ټول خصوصیات په کار دي د کومې نه چې دا راټول شوي او ډیزاین شوي وي - یو کارپس کښې د ژبې ټول اجزا لکه لفظونه، فقرې او جملې شاملې وي خو د یو کارپس سائیز د کارپس د لفظونو په شمېره انحصار کوي -
- په یو کارپس کښې د ژوند د هرې شعې نمائنده ګي په کار ده - لکه د سائنس، طب، ارتس، انجنیرنگ او د ټولو فنونونه اصطلاحات او لفظونه په کارپس کښې راغونډول په کار دي -
- پکار ده چې یو کارپس د چاپ شوو ټولو ذریعو نقل وي - په دې کښې د تلفظ، هجو، او املائي ټولو اصولو خیال ساتل په کار دي -
- یو کارپس په سافت فارم کښې محفوظول په کار دي چې د ژبې څیړن کاران دا ډېټا پېس د ژبې د تجزیې، د شماریاتي تجزیې او ترجمو ژباړو د پاره استعمالوي - [2]

۴ د کارپس ژبپوهنې یوه لنډه تاریخچه

کارپس جمع کول د ۱۹۶۰ء راهسې شروع شوي وو- په انګریزي ژبه کښې د تجرباتي ګرېمري څېړنو د پاره اولنې کارپس د Randolph Quirk Survey of English usage وه چې وروستو د ډیرو عشرو د پاره یو معیاري انګریزي ګرامر جوړ شو په ۱۹۸۵ء کښې دا Comprehensive Grammar of the English language په نوم چاپ شو او وروستو په ژبه کښې د تجرباتي څېړنو د پاره یو reference جوړ شو تر ډېره وخته دا کارپس کمپیوټري شومې نه وه - دا ډېر وروستو د ۱۹۸۰ء د عشرې په اخره کښې کمپیوټري شو- دا اوس د International Corpus

of English يعني ICE په نامه پېژندلې شي - د Quirk سروې د وېناوي اولیکلې ژبې ګډون وۀ - په یو ملین لفظونو کښې د وېناوي انګریزي پینځه لاکهه ټکي شامل وو - [8]

د ۱۹۷۰ء اخره کښې Jan Svartvik د وېناوي انګریزي کارپس کمپیوټر کښې واچوؤ - چې د London Lund په نوم پېژندلې شي - په دې کښې فونو لوجیکل او فونېټک معلومات هم ورکړي شوي دي - په ۱۹۶۰ء کښې دویم کارپس پراجیکټ د براؤن کارپس دے چې د Rhode جزیرې په Providence ښار کښې د براؤن یونیورسټي پروفېسر نیلسن فرانسز او هینري کاؤسیرا compile کړو - دا کارپس په یو ملین لفظونو مشتمل وۀ چې د امریکې پینځلس قسمونو پینځه سوه متنونو باندې مشتمل وۀ - دا کارپس ډېر احتیاط سره په منظمه طریقه را ټول شوی وۀ او د پروف ریډنگ غلطی ترې هم ویستلې شوې وې - [9]

دې پسې بل کمپیوټري کارپس په یو ملین لفظونو مشتمل لینکاسټراوسلو - برګن کارپس Lancaster/oslo-Bergen corpus وۀ چې په LOB پېژندلې شي - د ۱۹۶۵ء د کال چاپ شوي مواد د پینځه سوه متنونو نه دوه زره لفظونه راجمع کړي شوي وو - دا کارپس د Lancaster یونیورسټي Prof. Geoffrey Leech د نګرانی لاندې راټول شوی وۀ - دا په اصل کښې د امریکې د براؤن کارپس په مقابله کښې د برطانوي انګریزي د پاره راټول شوی وۀ - دا کارپس په دې وېب سائټ باندې سرچ کولې شي -

<http://khrhit.hit.uib.no/icame/manuals/lob/index.htm>[10]

دريم غټې کارپس پراجیکټ د انګریزي English Lexical Studies وۀ چې Edinburgh ښار کښې په ۱۹۶۳ء کښې شروع شو او په برمنګهم کښې ختم شو - ددې کارپس Principal Investigator ، John Sinclair وۀ ، دے اولنې ماهر وۀ چې کارپس ئې د Lexical Investigation د پاره استعمال کړو - دا پراجیکټ د وېناوي اولیکلي متن تقریباً په یو ملین لفظونو مشتمل وۀ د John Sinclair اخرنې رپورټ په ۱۹۷۰ء کښې د ټاټپ لیکلے یو خو کاپي تقسیم شوې - [11]

د ډکشنري جوړولو د پاره غټې غټې کارپورا هم جوړې شوې دي په ۱۹۸۰ء کښې Collins Publishers او برمنګهم یونیورسټي یو مشترکه پراجیکټ وۀ چې د پروفېسر John Sinclair د نګرانی لاندې جوړ شو چې د COBUILD په نامه پېژندلے شي او اوس تقریباً پینځه سوه ملینه لفظونو باندې مشتمل دے او په <http://www.mycobuild.com/about-collins-corpus.aspx> وېب سائټ باندې سرچ کولې شي -

دغه شان اکسفورډ یونیورسټي پېرس د خپلې ډکشنري د مرتب کولو د پاره اکسفورډ انګلش کارپس جوړ کړے دے چې تقریباً په دوه بلینه لفظونو مشتمل دے -

د ۱۹۹۰ء په شروع کښې برتښ نیشنل کارپس جوړ شو. دا په سلو ملینو لفظونو مشتمل وه چې نوي في صده ئې په لیکلي متن مشتمل دے او لس في صده ئې په وېناوي خبرواترو مشتمل دے. دا د ژبې د ډېټا یوه ښه ذریعه ده چې په انټرنیټ په www.natcorp.ox.ac.uk وېب سائټ باندې سرچ کېدے شي - [10]

۵ د کارپس قسمونه

د کارپس لاندیني قسمونه دي -

۱.۵ Reference Corpus

ریفرنس کارپورا د گڼو مقصدونو د پاره استعمالولے شي - په دې کښې د یوې ژبې معیاري لغات شامل وي - دا کارپس د ژبپوهانو د پاره په معنو پوهېدو یوه اهمه ذریعه ده - دې ته General corpora هم وئیلے شي - لکه برتښ نیشنل کارپورا BNC یا Bank of English.

۲.۵ Monitor Corpora

مانیټر کارپس یو داسې کارپس دے چې ژبه مانیټر کوي - دا په ژبه کښې بدلونو نه مانیټر کوي - دا وخت تروخته اپ ډېټ کولے شي - دا کارپس د لغاتو په بدلون کښې دلچسپي لري - لکه په ژبه کښې نوي لغات رانوتل او زاړه لغات نوې معنو کښې استعمالېدل یا د لفظونو، متلونو او اقوالو بدلون یا دوي کښې معنوي بدلون - دا قسم کارپورا زیات تر د اخبارونو او رسالو په ژبه مشتمل وي - د اخبارونو د ژبې خپل ځانگړے ستائل وي - د اخبارونو نه علاوه په دې کارپورا کښې کال په کال د ادب، تاریخ او سائنس نوي چاپ شوي کتابونه هم شاملولے شي - [12]

۳.۵ Parallel Corpora

Parallel کارپس ته ټرانسلیشن کارپس هم وئیلے شي - دا کارپس د یوې ژبې په اصلي متنونو مشتمل وي چې هم ددې متن په بله ژبه کښې ترجمه هم وي - کله کله دا ترجمه په ډېرو ژبو کښې هم وي -

Historical Corpora ۴.۵

تاريخي کارپورا په تاريخي موادو مشتمله وي - لکه دمثال په توگه Helsinki کارپس يا Archer -

Regional Corpora ۵.۵

دا د يوې علاقې د ژبې په لغاتو مشتمل وي - لکه WCNZE يعني Wellington Corpus .of written new Zealand English

Learner Corpora ۶.۵

د يوې ژبې د زده کړې د پاره چې غېر مقامي خلقو ته کوم ضروري لغات زده کول په کار وي - دا قسم کارپورا په ټولو هغه لغاتو مشتمله وي لکه ICLE ، International Corpus .of Learner English

Multilingual Corpora ۷.۵

دې قسم کارپورا کښې د گڼو ژبو او کم از کم د دوو ژبو لغات شامل وي - دې کښې اکثر د پرتليزه څېړنو د پاره يو متن په دوو يا ډېرو ژبو کښې ورکړې شوي وي -

Spoken Corpora ۸.۵

عام کارپورا کښې وېناوي او ليکلي متنونه دواړه شامل وي - خو په Spoken کارپورا کښې صرف وېناوي مواد شامل وي - لکه London-Lund Corpus of Spoken English

Internet Corpora ۹.۵

نن سبا په ژبپوهانو کښې دا Trend جوړ شوي دي چې انټرنېټ هم د کارپورا په توگه استعمالوي - دا په خصوصي توگه هغه وخت فائده مند وي چې کله مونږ د يو لفظ په حقله کتل غواړو چې ايا دا لفظ شته او کنه او په کومو کومو مطلبونو کښې استعمالېږي -

Annotated Corpora ۱۰.۵

په Annotated کارپورا کښې څه ژبنی تجزيې د منځکښې نه شوې وي، لکه د لفظونو وېش يا د جملو تجزيه -

لنډه دا چې کارپس ژبپوهنې ډېره وده موندې ده او نن سبا په ژبپوهنه کښې په گڼو اړخونو کښې استعمالېږي - په کارپس ژبپوهنه کښې نړيوالې تحقيقي مجلې چاپ کېدل هم شروع شوي دي لکه Corpus Linguistics & Linguistic theory International Journal of Linguistics، Corpora Linguistics، Corpus studies in corpus او ددې نه علاوه د کتابونو ډېر سيريز د کارپس ژبپوهنې اوځي چې په کښې د ژبې او کمپوټر په تعلق بحث شوي وي لکه External links او English corpus linguistics ددې نه علاوه څه هم دې چې مختلف کارپورا پرې مطالعه کېدې شي، لکه Penn parsed corpora of historical English Linguistic Data consortium او corpora discussion list وغېره -

حوالې

- [۱] http://www.as.uni.heidelberg.de/personen/nesselhauf/files/corpus-linguistics_practical_introduction.pdf (accessed June 2, 2017)
- [۲] <http://www.eolss.net/sawple-chapter/co4/e6-91-17> PDF (accessed May 19, 2017)
- [۳] Tony McEnery et al., *Corpus-Based Language Studies. An Advance Resource Book.* (London & New York: Routledge Taylor & Francis Group 2006) P. 4
- [۴] Ibid- P-3
- [۵] Michael Stubbs, "Language Corpora" in *The Handbook of Applied Linguistics*, Eds: Alan Davies & Catherine Elder (Malden USA: Blackwell Publishers, 2004) P. 107
- [۶] Ibid- P. 124
- [۷] Randi Reppen & Rita Simpson-Vlach, "Corpus Linguistics" in *An Introduction to Applied Linguistics* (UK: Hodder Education An Hachette UK Company, 2011) P. 95
- [۸] Wolfgang Teubert & Anna Cermakora, *Corpus Linguistics A short Introduction* (London: Continuum Publishers, Reprint 2008) P. 51
- [۹] Ibid P- 52
- [۱۰] Howard Jackson & Peter Stockwell, *An Introduction to the Nature & Function of Language* (London & New York: Continuum International Publishing Group 2010) P. 12
- [۱۱] Wolfgang Teubert & Anna Cermakora, *Corpus Linguistics A short Introduction* (London: Continuum Publishers, Reprint 2008) P. 54
- [۱۲] Ibid, P. 72