

Bilgisayar Ortamında Bir Derlem Geliştirme Çalışması

Bilge Say,
Enformatik Enstitüsü
Bilişsel Bilimler Ana Bilim Dalı
Orta Doğu Teknik Üniversitesi
bsay@ii.metu.edu.tr

Umut Özge,
Enformatik Enstitüsü
Bilişsel Bilimler Ana Bilim Dalı
Orta Doğu Teknik Üniversitesi
umut@ii.metu.edu.tr

Kemal Oflazer,
Mühendislik ve Doğa Bilimleri Fakültesi
Sabancı Üniversitesi
oflazer@sabanciuniv.edu

Öz: Bu bildiri dilbilim ve bilgisayarlı doğal dil işleme çalışmalarına kaynak olmak üzere elektronik ortama geçirilen günümüz Türkçesini yansıtan metin örneklerinin işaretlenmesiyle oluşturulan bir derlemin geliştirilme süreci özetlenmektedir. Özellikle, derlemin tasarımı ve kullanımı açısından gerçekleştirilen bilişim süreçlerinin detayları derlem oluşturma sürecinin genel perspektifi içinde verilecektir.

Ana Konu: Araştırmada bilişim.

Anahtar Sözcükler: Derlem, bilgisayarlı dilbilim, TEI, XML, XCES.

1. Akademik Bir Kaynak Olarak Derlem

Derlem (bütünce, corpus (İng.)) belli prensipler çerçevesinde özel veya genel amaçlı metin ya da konuşma parça ya da bütünlerinin, üzerinde yapılacak araştırmaya uygun işaretlemelerle beraber bir araya getirilmesinden oluşan bütündür. Günümüz derlemlerinin elektronik ortamda tutularak, erişim ve kullanım kolaylığı sağlanması yaygındır (Kennedy,1998). Bu da derlem oluşturma çabalarının bir akademik bilişim aktivitesi olarak ele alınmasını gerektirir.

Bir derlemden nasıl yararlanılabileceğine kısaca değinmek bu emek yoğun sürecin değerlendirilmesinde yararlı olacaktır. Dilbilimin pek çok altsahasında bir derlemden yararlanmak bilimsel bir metod olarak araştırmalara katkıda bulunmaktadır: bir sözlük oluşturulmasında günümüzün yeni sözcüklerinin nasıl bağlamlarda kullanıldığını araştırmak; sözdizimsel bir kuramın savlarını yaşayan dilin kesitlerinden örneklerle ve istatistiklerle desteklemek; bir dil öğrenimi sınıfında dilin yapılarına ve sözcüklerine örnekler vermek bu kullanımların sadece bazılarıdır (McEnery ve diğerleri, 1996). Doğal dil işleme alanında son yıllarda ağırlıklı olarak kullanılan istatistikî modellerin başarılı kullanımı için bol miktarda veriye, yani bir derleme ihtiyaç duyulur (Manning ve Schütze, 1999).

Özellikle Batı dillerinde son on-onbeş yılda tüm araştırmacıların kullanımına açık derlem

oluřturma alıřmaları meyvelerini vermiř; İngiliz Ulusal Derlemi (BNC, 2001), ek Ulusal Derlemi (Cermák, 1997) gibi 100 milyon szckle hem szel hem yazılı dili temsil eden kapsamlı derlemler oluřturulmuřtur. Her ne kadar Trke zerine alıřan dilbilimciler arasında yaptığımız bir ankette, katılanların %41'i Trke bir derlem kullandıklarını belirtmiř olsalar da bu derlemlerin ođu kiřisel olarak biraraya getirilmiř, kek lekli, elektronik olmayan ve tm akademisyenlerin ulařımına aık olmayan derlemlerdir. Oysa hem Trke dilbilim alıřmalarının ve hem de Trke zerine yapılan dođal dil iřleme alıřmalarının (Oflazer 1997) geliřimi deđiřik trleri ieren, elektronik ortamda, telif haklarına saygılı, tm arařtırmacılara aık ve bilgisayar ortamında sorgu yazılımlarıyla desteklenen birden ok derleme ihtiya duyulduđunu gstermektedir.

2. ODT Trke Derlemi Tasarımı

Bu blmde hem ODT Trke derlemi tasarımının genel đeleri, hem de bu derlemin bir altderlemi olan ODT-Sabancı ađa yapılı derleminin genel hatları verilecektir.

2.1 Genel Tasarım đeleri

Bu gereksinmeden yola ıkan ODT Trke Derlemi Geliřtirme projesi kısıtlı bte ve personel olanaklarıyla nceki blmde bahsedilen gereksinimlere belirli sınırlamalar erevesinde yanıt olacaktır. Bu sınırları kısaca gzden geirelim: Derlem sadece yazılı dili iermekte, daha yođun emek gerektirdiđinden szl dil kapsam dıřı tutulmaktadır. Trk dilinin geliřiminden ok gnmzdeki durumunu temsil etmeyi amalayarak 1990 sonrasında yayınlanan eserlerden rneklemler alınmaktadır. Derlem dilin deđiřik ynlerinin incelenmesi aısından deđiřik trlerde (Bkz. Blm 3) rneklemler iermekte, ancak bu trlerin derlem ierisindeki yzdesi dilin retimi ve tktimine dair bir arařtırmaya dayanmamaktadır (Biber, 1993). Byklk aısından 2 milyon szck hedeflenmekte; bu Brown derlemi gibi 80'lerde sık kullanılan 1 milyon szcklk derlemlere gre byk, ancak szlkbilimsel alıřmalar iin alt sınır kabul edilen 10 milyon szck gereksinimine gre kek bir miktardır (Kennedy, 1998). Burada da sınırlayıcı etken eldeki iř gc olanađının kısıtlılıđı olmuřtur. Bir rneklemdaki szck sayısı Brown ve benzerlerinde olduđu gibi 2000'dir ki bu seimde de belirleyici etken, yayınevlerinin genelde uzun rneklemlere telif hakları aısından sıcak bakmaması olmuřtur. Derlemin bir kısmı biimbirimsel ve szdizimsel iřaretlenmek zere ayrıca projelendirilmiřtir (Bkz. Blm 2.2). Derlemin tm dnyada metin iřaretleme standardı olarak kabul gren TEI (Text Encoding Initiative) standardının derlemlere zel bir XML uygulaması olan XCES'le iřaretlenmektedir (Bkz. Blm 3).¹

Bu tasarım parametrelerinin belirlenmesiyle bařlayan derlem oluřturma alıřmalarında derlem oluřturma srecinde iki ana biliřimsel aktivite ortaya ıkmıřtır: Metin rneklemlerinin elektronik ortama alınması ve iřaretlenmesi ile dilbilimcilerin kolay kullanımı iin genel amalı bir sorgu yazılımı geliřtirilmesi. Blm 3'de bu sreler detaylı olarak tanıtılacaktır.

2.2 ODT-Sabancı Ađa Yapılı Derlemi Tasarımı

Penn Ađa Yapılı Derlemi (Penn) gibi derlemler son yıllarda teorik dilbilim ve dođal dil iřleme alıřmalarında, geliřtirme ve geerliliđi lme aısından nemli kaynaklar olmuřlardır. Bu aıdan Trke derlemin bir kısmının biimbirimsel ve szdizimsel olarak iřaretlenmesi Sabancı niversitesi ile iřbirliđi iinde ayrı bir proje olarak gerekleřtirilmektedir. Trkenin biimbirimsel (morfolojik) olarak zengin bir

¹ Derlem tasarımı ile ilgili detaylı bilgi iin bkz. Atkins ve diđerleri (1993).

dil olması, tümcenin öğeleri arasında sözcüklere eklenen sonlu bir işaret kümesine dayanan değil, “çekimsel grup” adını verdiğimiz biçimbirimsel altyapılara dayanan bağılılık (dependency (İng.)) yapılarının işaretlenmesinin dilbilimsel ve bilişimsel açıdan anlamlı bir gösterim olduğu sonucunu getirmiştir. (Oflaz ve diğerleri, 2000). Bu sonuçtan yola çıkan bir tasarımla şu ana kadar bir işaretleme kılavuzu ve işaretlemenin kolaylıkla gerçekleştirilmesi için bir yazılım geliştirilmiştir. Bundan sonra biçimbirimsel olarak ayrıştırılan öğelerin çoklu yapılarının indirgenmesi (disambiguation (İng)), öbek yapısı gösteren “şırıl şırıl” gibi dizilerin işaretlenmesi ve kalan öğeler arasında özne, nesne, belirteç gibi fonksiyonel bağların yarı otomatik işaretlenmesi planlanmaktadır.

3. Derlem Yapım Süreci

Bu bölümde derlemin yapım aşamasında yer alan süreçlere değineceğiz. Derlem yapımı, *metin toplama*, *işaretleme* ve *kontrol* olmak üzere üç ana süreci kapsamaktadır. Takip eden bölümlerde bu süreçleri ayrı ayrı ele alacak, süreçlerin detayları, uygulanan standartlar ve kullanılan yazılımlar hakkında bilgi vermeye çalışacak, bu kısmın son bölümünde de derlemi sorgulamak amacıyla yönelik, henüz yapım aşamasındaki derlem sorgu yazılımının özelliklerine değineceğiz.

3.1 Metin Toplama

Metin toplama süreci, derleme alınacak eserlerin telif hakkı sahiplerinden yazılı izin alınması ile başlayıp, metinlerin bir sonraki süreç olan işaretleme süreci öncesi elektronik ortamda "unicode salt metin" formatında hazır bulunmasıyla sona erer.

Derleme alınacak metinlerin belirlenmesinde rol oynayan faktörlerden biri, derlem içeriğinin türlere göre dağılımının dengeli olmasıdır. Bu dağılımı takip edebilmek amacıyla 14 metin türü kategorisi belirlenmiştir. Bu makalenin yazıldığı sırada 1.000.000 sözcüğe (2000'er kelime) 250 örneklem) ulaşan derlem içeriğinin türlere göre dağılımı Tablo 1'deki gibidir.

Tür	Dağılım (%)
Roman	24
Öykü	21
Makale	16
Deneme	14
Araştırma-İnceleme	12
Gezi	4
Söyleşi	2
Diğer	7

Tablo 1: Derlem içeriğinin türlere göre dağılımı.

Yukarıdaki tablodaki “Diğer” başlığı *köşe yazısı*, *referans*, *anı*, *yaşam öyküsü*, *özyaşam öyküsü*, *kişisel gelişim* ve *ders kitabı* türlerini kapsamaktadır.

Derleme dahil edilecek eserler, tercihen elektronik ortamda yayıncıdan sağlanmaya çalışılmıştır. Bunun mümkün olmadığı durumlarda üniversite kütüphanelerinden veya yayıncı kuruluştan temin edilmektedir. Bu eserlerin içinden seçilen örneklem HP ScanJet 6200C tarayıcı ile taranıp, OCR tekniğiyle "salt metin" formatında elektronik ortama aktarılmaktadır. Tarama işlemi sırasında bazı

karakterlerin yanlış tanınması sonucu oluşan hatalar elle düzeltilmekte, bu yolla elektronik ortama aktarma işlemi sırasında orijinal metne sadık kalınmaktadır. 2000 kelimededen oluşan bir örneklemin kaynağından taranarak elektronik ortama aktarılması bir çalışanın 1 ½ saatini almaktadır.

3.2 İşaretleme

Derlemimizde bulunan elektronik metinlerin kodlanmasında, metinlerin basılı buldukları kaynaktan bağımsız olarak paragraf, tırnak, listeleme ve benzeri öğelerinin ve ayrıca künye bilgilerinin kodlanması için oluşturulan TEI (Text Encoding Initiative) uygulaması olan XCES'i kullanmaktayız.

1980'lerin sonuna doğru internet'in de yaygınlaşmaya başlamasıyla, sayıca artan elektronik metinlerde bulunan bilginin, donanım, yazılım ve uygulamalardan bağımsız bir şekilde paylaşılması ve verimli bir şekilde işlenilmesi sorunu doğmuştur. Bu sorunu çözmek için bir grup araştırmacının başlattığı 7 yıl süren çalışma 1994 yılında ilk resmi TEI kılavuzunun yayınlanmasıyla meyvesini vermiştir (Sperberg-McQueen ve Burnard, 1994). TEI, metinlerin elektronik ortamda gösterimine evrensel bir standart getirmeyi amaç edinmiş uluslararası akademik bir araştırma hareketidir.

TEI bizlere elektronik metinlerin kodlanmasına yönelik çok temel düzeyde SGML işaretlerinden ve ortaya çıkacak SGML dokümanlarının düzgün yapılandırılmış olma koşullarını belirleyen DTD (Document Type Definition)'lerden oluşan bir standart sunmaktadır. Bu standardın eldeki amaca uygun şekilde özelleştirilerek geliştirilmesi gerekmektedir. Dilbilimsel derlemeler için TEI kılavuzları ile uyumlu bir SGML (Standard Generalized Markup Language) uygulaması olan ve dilbilimsel derlemelerin genel mimarisini ve dilbilimsel işaretleme standartlarını belirleyen CES (Corpus Encoding Standard) elektronik metin kodlama kılavuzu geliştirilmiştir (Ide, 1996).

Şu ana kadar bahsettiğimiz metin kodlama standartlarının hepsi SGML uygulamalarıdır. Günümüzde web üzerindeki veriyi görüntülemekte neredeyse evrensel bir dil olarak kabul görmüş HTML'in de tabanını oluşturan SGML, ilk başta yukarıda bahsettiğimiz standardizasyon için en uygun dil olarak görülüyordu. Fakat çok geçmeden SGML'in bir takım eksileri göze çarpmaya başladı. Bu eksiler: (i) SGML'in son derece karmaşık bir sözdiziminin olması; (ii) SGML dokümanlarının DTD olmadan ayrıştırılamaması; (iii) SGML dokümanlarının kısmen ayrıştırılma şansının olmaması ve bunun özellikle büyük dokümanlarda zaman kaybına yol açması olarak özetlenebilir (DeRose, 1999). Bu noktadan hareketle, yukarıda özetlediğimiz eksilerden arındırılmış, SGML'in gücü ve esnekliği ile HTML'in basitliğini birleştiren XML dili geliştirilmiştir (Bray ve diğerleri, 1998).

Biz de bu bağlamda CES'in, bilişim ve internet teknolojisinde hızla yaygınlaşan XML diline adapte edilmiş hali olan XCES'i derlemimizde işaretleme standardı olarak kullanmayı uygun bulduk. CES'in XML'e adapte edilmesi halen devam eden bir süreçtir. (Welty ve Ide, 1999, s.62).

Derlemi oluşturan ve her biri 2000 sözcüklük bir metin ve bu metni işaretlemede kullanılan XCES işaretlerini içeren bir XML dokümanı olan örneklem, *başlık* (header) ve *gövde* (body) olmak üzere iki bölümden oluşmaktadır.

Başlık bölümünde örneklemin alındığı kaynağın detaylı bibliyografik bilgilerinin yanı sıra, örneklemin yazılı olduğu dosyanın adı, dosyanın büyüklüğü, örneklemin içerdiği sözcük sayısı (işaretler hariç), kontrol işlemi (Bkz. Bölüm 3.3) sırasında eğer metnin orijinali üzerinde bir değişiklik yapıldıysa, ne tür bir değişikliğin, kim tarafından ne zaman yapıldığı gibi bilgiler yer almaktadır. Şekil 1'de tipik bir örneklem başlığı görülmektedir.

```
<cesHeader>
<fileDesc>
  <titleStmnt>
    <h.title>00017113</h.title>
  </titleStmnt>
  <extent>
    <wordCount>2008</wordCount>
    <byteCount>17929</byteCount>
  </extent>
```

```

<sourceDesc>
<biblStruct>
<analytic>
<h.title>Anadolu Dağlarının 'Bitki Avcısı': Prof. Dr. Turhan BAYTOP</h.title>
<h.author>Nalân MAHSERECİ</h.author>
</analytic>
<monogr>
<h.title></h.title>
<h.author></h.author>
<edition></edition>
<imprint>
<publisher>Bilim ve Ütopya</publisher>
<pubDate>Mart 2000</pubDate>
<pubPlace>İstanbul</pubPlace>
</imprint>
<idno>1301 - 6717</idno>
<biblScope>69</biblScope>
</monogr>
</biblStruct>
</sourceDesc>
</fileDesc>
<profileDesc>
<textClass>
<catRef>Makale</catRef>
</textClass>
</profileDesc>
<revisionDesc>
<change>
<changeDate>12.10.2000</changeDate>
<respname>Sedef</respname>
<h.item>The header part was changed.</h.item>
</change>
</revisionDesc>
</cesHeader>

```

Şekil 1: Tipik bir örneklem başlığı.

Örneklemin *gövde* bölümünde üç aşamalı bir standart olan CES'in asgari işaretlemeyi belirleyen birinci aşaması uygulanmıştır (Ide, 1996, Bölüm 4.1). *Gövde* kısmında kullandığımız işaretlere kısaca göz atalım.

(a) Üst paragraf düzeyi işaretlerinin bazıları:

<text> metinleri işaretler.

<body> metin içinde bütünlüğü olan parçaları işaretler. (örneklemin ana kısımlarında biri olan *gövde* (body) kısmıyla karıştırılmamalıdır.)

<opener> metinlerin başlangıcındaki, tarih, anahtar sözcükler ve benzerlerini işaretler.

<head> metin, liste şiir gibi yapıların başlıklarını işaretler.

(b) Paragraf Düzeyi İşaretlerinin bazıları:

<p> paragrafları işaretler.

<q> tırnak içine alınmış kısımları işaretler.

<hi> normal karakter formatı dışında yatık, koyu, altıçizili gibi vurgulanmış sözcükve sözcük öbeklerini işaretler.

<poem> şiirleri işaretler.

<table> tabloları işaretler.

<list> listeleri işaretler.

<note> metin içinde geçen her türlü notu işaretler.

<abbr> kısaltmaları işaretler.

<date> tarihleri işaretler

Şekil 2'de Nihal Yeğınobalı'nın Can Yayınları tarafından basılmış "Sitem" adlı romanından alınmış bir örneklemin *gövde* kısmının kısaltılmış hali görülmektedir.

```
<text>
<body>
.
.
<p>Oktay biraz önce, <q>Hadi biz de Sitem'in yanına gidelim,</q> demişti. Sitem'in, kucağında Tomurcuk Beyle Yılanlı İncirlerden yana gittiğini o da görmüştü çünkü. Ben omuz silkmekle yetindim, Oktay da üstelemedi. Sitem ikimizin yüzüne karşı da görünmez kapılar kapamıştı. Benim de elinden kayıp gidivermemden korkan Oktay beni <hi>oyalamak</hi> için geçen yaz Giray Ağabeysiyle Kirazlı Yaylaya yaptıkları bir gezintiyi anlatmaya başladı.</p>
.
.
<p>O gün ve sonrasında olanları elbet sana da anlatmışlardır, Dalya. Gene de o kargaşa, o şaşkınlık, o panik, o kafa karmaşası yaşanmadan bilinemez...</p>
.
.
</body>
</text>
```

Şekil 2: Tipik bir örnekleme gövdesi.

İşaretleme süreci, elektronik ortama aktarılmış ham metnin, grubumuzca C'de yazılmış, en sık kullanılan “<p>” ve “<q>” işaretlerini otomatik olarak metne yerleştiren bir yazılımla işlenmesiyle başlar. Bu safhadan sonraki işaretleme işlemlerinde yine grubumuzca Borland C++ Builder'da yazılmış XCESEdit adını verdiğimiz bir XML editörü kullanılmaktadır. Grafik kullanıcı arayüzüne sahip bu program örneklemin *gövde* kısmında işaretleme sürecinin imlecin bulunduğu yerlere sağ fare tuşunu kullanarak gerekli işaretleri koyabilmesine olanak sağlamaktadır. İşaretleme süreci hiçbir işareti kendisi yazmadığı için işaretleri yanlış yazması olasılığı ortadan kalkmaktadır. *Gövde* kısmının işaretleme sürecinin tamamlanmasının ardından, editörün *başlık* oluşturma işlevi kullanılır. *Başlık* oluşturma işleminde editör işaretleme sürecinin önüne, alanları Şekil 1'deki gibi bir *başlığın* işaretlerini kapsayan bir form getirir. Bu form doldurulurken işaretleme süreci herhangi bir alana uygun olmayan bir bilgi girerse yazılım tarafından uyarılır. (Ör. örnekleme sürecini oluşturan metnin türünün belirten <catRef> alanına projede kullanılan türler dışında herhangi birşey girerse yazılım hata mesajı verir). Form doldurulduktan sonra örneklemin *başlığı* otomatik olarak oluşturulur. Tüm bu işlemler bittikten sonra yine aynı editör kullanılarak örnekleme süreci ayrıştırılır. Editör, açılan işaretlerin kapanıp kapanmadığını, işaretlerin içiçe geçip geçmediğini kontrol ederek hata olan yerleri işaretleme sürecine bildirir. Böylece düzgün yapılandırılmış bir XML dokümanı oluşturulmuş olur. Bahsedilen yazılımlar kullanılmaya başlanmadan önce bir işaretleme sürecinin yaklaşık 1 saat 15 dakikasını alan işaretleme süreci yazılımların kullanılmasıyla 1/2 saatte tamamlanabilmektedir.

3.3 Kontrol

Yazılım desteğine rağmen gözden kaçan bir takım hatalar olabilmektedir. Bu nedenle bir önceki bölümde anlattığımız şekilde oluşturulan örneklemler son bir defa daha kontrol edilmektedir. Kontrolü yapan kişi bu safhada örnekleme sürecini orijinal metinle karşılaştırır. Orijinal metinde imla hatası yapıldığından şüphelendiği durumlarda Türk Dil Kurumu'nun 2000 yılı baskısı "İmla ve Yazım Kılavuzu"na başvurarak gerekli düzeltmeyi yapar ve bu düzeltmeyi de metin içinde kodlar.

Her örnekleme, aynı zamanda örneklemin yazılı olduğu dosyanın da ismi olan 8 haneli bir kod numarası verilmektedir. Kod numarasının ilk 5 hanesi örneklemin alındığı kaynağı belirtmektedir. Kaynaklar derleme dahil edilme sıralarına göre 00001'den başlayarak numaralandırılmıştır. İlk beş haneyi takip eden hane ise o örneklemin eldeki kaynaktan alınan kaçınıcı örnekleme sürecini belirtmektedir. Bir sonraki hane örnekleme sürecini hangi çalışanın işaretleme sürecini göstermektedir. (Proje çalışanları tek basamaklı bir

sayıyla kimliklendirilmiştir.) Sekizinci ve son hane ise örneklemin kontrolünü yapan kişiyi belirtmeye ayrılmıştır. Örneklemlerin kodlanmasını bir örnek ile açıklayacak olursak; 00125273 kodlu bir örneklem “125” no’lu kaynaktan alınan ikinci örneklem olup “7” kimlik numaralı çalışan tarafından işaretlenip “3” kimlik numaralı işaretleyici tarafından kontrol edilmiştir. Bu kodlama sistemi ve *başlık* kısmında kodlanan değişiklik bilgileri sayesinde hataların kimler tarafından yapıldığı takip edilebilmekte, çalışanların yaptıkları hatalar konusunda uyarılmasıyla derlemin kalitesi yüksek düzeyde tutulmaktadır.

3.4 Derlem Sorgu Yazılımı Tasarımı

Derlem, kullanıcılara Java dilinde yazılmış, kullanıcıya grafik bir arayüz kullanarak derlemi sorgulama imkanı sağlayacak bir derlem sorgu yazılımı ile birlikte dağıtılacaktır. Bu bölümde henüz yapım aşamasında olan bu yazılımın bazı özelliklerine değineceğiz.

Derlem sorgu yazılımının geliştirilmesinde göz önünde bulundurduğumuz iki ana tasarım kriteri şunlardır;

- (i) Kullanıcıların derlemden elde etmeyi amaçladıkları bilgi türündeki çeşitliliği hesaba katarak, mümkün olduğu kadar geniş bir yelpazede sorgu yapabilme imkanı tanımak.
- (ii) Sorgular karşılığında derlemden çıkarılan bilginin kolay analiz edilmesini sağlayacak, bir taraftan başlangıç düzeyindeki bilgisayar kullanıcılarının bile zorluk çekmeden kullanabileceği kadar basit, diğer taraftan da ileri düzey kullanıcılara kullanım esnekliği sağlayan bir grafik arayüz sunabilmek.

Derlem sorgu yazılımı istemci/sunucu mimarisinde, nesne tabanlı yaklaşımla Java programlama dilinde yazılmaktadır. Java dilinin seçilmesinin nedeni, sözkonusu programlama dilinin değişik platformlarda büyük ölçüde sorunsuz çalışabilen, unicode karakter kümesi kullanması itibariyle Türkçe karakterlerde sorun çıkarmayan ve web üzerinde çalışabilen uygulamalar yaratılmasına olanak sağlamasıdır.

Kullanıcıların derlemden sorgulamak ihtiyacını duyacağını düşündüğümüz ve sorgulama yazılımının olanak vereceği sorgular şöyledir;

1. derleme göz atma;
2. sözcük sorgusu;
3. XCES işaretlerine göre sorgulama;
4. bağlılık yapısına göre sorgulama (ağaç yapılı derlem);
5. biçimbirimsel özelliklere göre sorgulama;
6. düz terim (regular expression(İng.)) sorgulama;
7. çapraz sorgular yapabilme.

Kullanıcı, bu sorgular karşısında elde ettiği sonuçları tarih ve saatiyle birlikte kaydedebilecektir. Kullanıcı, yazar adı, yayınevi, tür, basım yılı vb. kriterler belirterek yalnızca bu kriterlere uyan örneklemi kapsayan bir alt derlem oluşturabilecek, ve sorgularını bu alt derlem üzerinde yapabilecektir. (ör. 1995-1996 yılları arasında basılmış, öykü türündeki metinlerdeki alıntılar sorgulama.) Ayrıca kullanıcı bu alt derlemleri kaydedebilme, dolayısıyla istediği zaman tekrar kullanabilme imkanına da sahip olacaktır. Yazılım, prototip olarak tamamlandığında bir grup dilbilimci tarafından denenecek ve gerekli düzeltmeler yapılacaktır.

4. Sonuç

Türk Dili açısından bir ilk olacak olan bu derlem projesi, Mayıs 2002’de ağaç yapılı altderlem kısmı dışında sona erecektir. Yaşadığımız deneyimin bu bildiriyle aktarılmasının daha kapsamlı ya da özel amaçlı başka derlemlerin gerçekleştirilmesine katkıda bulunacağını umuyoruz.

Teşekkür

Projede halen veya başlangıçta fikren ve/veya emeği ile katkısı olan öğretim üyeleri Prof. Dr. Wolf König, Prof. Dr. Deniz Zeyrek, Doç. Dr. Cem Bozşahin, Y. Doç. Dr. Ümit Deniz Turan, Y. Doç. Dr. Margaret Sönmez, Dr. Ayşenur Birtürk, Dr. Dilek Hakkani-Tür ve Dr. Gökhan Tür’e; araştırma görevlileri Barış Şükrü Demiral, Barış Çağrı Genç, ve Filiz Yılmaz Bican’a; ve “işaretleyicilerimiz” Sedef Akgül, Aygün Boduroğlu, Deniz Cantürk, Devrim Saran ve Barış Şara’ya; bizi maddi olarak destekleyen ODTÜ AFP (AFP No: 99-06-04-02) ve TÜBİTAK’a (EEEAG Proje No: 199E026); metin toplamada bize yardımcı olan tüm yayınevi ve kuruluşlara (Can Yayınları, İletişim Yayınevi, Bilgi Yayınevi, Kuraldışı Yayınevi, Adam Yayınları, İşbankası Kültür Yayınları, Yapı Kredi Yayınları, Bilim ve Ütopya, Doğu-Batı, Atlas ve Bütün Dünya dergileri.); ve Türkçemize “derlem” sözcüğünü kazandıran ve deneyimlerini bizimle paylaşan Prof. Dr. Aydın Köksal’a teşekkür ederiz.

Kaynakça

- Atkins, S, J. Clear, ve N. Ostler. 1993. *Corpus Design Criteria*. Literary and Linguistic Computing 8(4).
- Biber, D. 1993. *Representativeness in Corpus Design*. Literary and Linguistic Computing 8(4).
- BNC. British National Corpus. <http://info.ox.ac.uk/bnc/>, Aralık 2001 haliyle.
- Bray, T., J. Paoli, ve C. M. Sperberg-McQueen. 1998. *Extensible Markup Language (XML) 1.0. W3C Recommendation*, Şubat.
- DeRose, S.J. 1999. *XML and the TEI*. Computers and the Humanities 33, 11-30.
- Cermák, F. 1997. *Czech National Corpus: A Case in Many Contexts*. International Journal of Corpus Linguistics 2(2).
- Ide N. 1996. *Corpus Encoding Standard: Document CES 1, Sürüm 1.4*, Ekim. <http://www.cs.vassar.edu/CES/>
- Kennedy, G. 1998. *An Introduction to Corpus Linguistics*. Longman.
- Manning, C. D. ve Schütze. H. 1999 *Foundations of Statistical Natural Language Processing*. Cambridge MA: MIT Press.
- Mc Enery, T ve A. Wilson. 1996. *Corpus Linguistics*. Edinburgh: Edinburgh University Press.
- Oflazer, K. 1997. *Natural Language Processing Research in Turkey*. Proceedings of 3rd Telri European Seminar, Montecatini, Italy, Oct 16-18.
- Oflazer, K, B. Say, D. Z. Hakkani-Tür, G. Tür. 2000. Building a Turkish Treebank. Abeille A (haz.) *Building and Exploiting Syntactically Annotated Corpora*.
- Penn Treebank. <http://www ldc.upenn.edu/>, Aralık 2001 haliyle.

Sperberg-McQueen, C. M., ve L. Burnard, (haz.) 1994. *Guidelines for Electronic Text Encoding and Interchange (TEI P3)*, Chicago, Oxford: Text Encoding Initiative.

Welty, C., N. Ide. 1999. *Using the Right Tools: Enhancing Retrieval from Marked-up Documents*. *Computers and the Humanities* 33, 59-84.