

VERİ SIKIŞTIRMA PROJESİ

Osman Şimşek
Bilgisayar Mühendisliği
Kocaeli Üniversitesi
180202048
simsekosman15@hotmail.com

Yener Emin Elibol
Bilgisayar Mühendisliği
Kocaeli Üniversitesi
180202054
yenereminelibol@gmail.com

Veri Sıkıştırma Projesi karakter tabanlı sıkıştırma algoritması kullanarak kullanıcının girdiği verileri sıkıştırmayı sağlayan projedir.

I. GİRİŞ

Veri Sıkıştırma Projesi karakter tabanlı sıkıştırma algoritması olan Lz77 algoritmasını kullanarak kullanıcının proje dosyalarının konumunda bulunan "Metin.txt" dosyasına girdiği verileri sıkıştırır. Lz77 algoritmasında ki temel amaç metin içinde tekrar eden kısımları tespit edip token denilen yapılar kullanılarak tekrar etmeyen ve tekrar eden kısımlara tokenler koyup sıkıştırılma yapması beklenmektedir. Tokenlerin temel amaçlarından biride sıkıştırma yapılan dosyanın eski haline getirilmesi için kolaylık sağlamaktır.

II. TEMEL BİLGİLER

.Program C programlama dilinde gerçekleştirilmiş olup geliştirme ortamı olarak CodeBlocks kullanılmıştır.

III. TASARIM

A. Algoritma

Kullanıcı sisteme veri girişi yapıp program çalıştırıldığında sistem kullanıcının girdiği verileri Metin.txt dosyasından sisteme çekerek Lz77 algoritması uygulamak üzere fonksiyona göndermektedir. Lz77 algoritmasının amacı, metin belgelerindeki tekrar eden bölümleri ortadan kaldırmak ve buna yönelik metin belgesini tokenler yardımıyla sıkıştırmaktır. Lz77 metnin içerisinde Lz77 yaklaşımında sözlük, daha önce kodlanmış serinin bir parçasıdır. Algoritmadaki arama tamponunun büyüklüğü, daha önce kodlanmış serinin ne büyüklükte bir parçasında arama yapılacağını belirler. Arama tamponu büyütüldükçe, sıkıştırma oranı artar, fakat aynı zamanda sıkıştırma zamanı da artar. Arama tamponu bu programda en fazla 4096 karakter öncesinden arama yapmaya başlayabilmektedir ve ileri tampon miktarı ise en fazla 16 karakter uzunluğundadır.

Arama tamponu ile ileri tampon arasında benzer bir kelime bulunursa bu benzerliğin ne kadar uzunlukta olduğu bulunur. Bulunan benzerliğin ana konumdan ne kadar gerisinde başladığı bilgileri ileri tampon ve arama tamponu yardımı ile bulunur. Bulunan bu değerleri token oluşturulup içerisinde bulunan 16 bit boyutundaki bir değişkene aktarılmaktadır. Bunlardan 12 si arama tamponuna 4 ü ise ileri tampon bilgisini tutmaktadır. Eğerki arama ve ileri tampondan gelen değerleri int tipinde tutsaydık tokenimizin boyutu artıcağı ve dosyamızın boyutuda artıcağı bunun yerinde 16 bitlik bir değişken kullanarak sıkıştırma yapılmaktadır. Tokenlerimizin içerisinde eşleşmeden sonra gelen ilk değer tutulmaktadır. [10,7,a] Tokenlerin genel gösterimleri bu şekildedir ilk değer arama tampon bilgisi

ikinci değer ileri tampon bilgisi 3. değer eşleştirmeden sonra gelen ilk karakterdir eğerki eşleştirme bulunamadıysa o anki karakterin tokeni [0,0,a] gibidir. Her adımda token oluşturulmaktadır eşleşme yakalansa da yakalanmasa da. Eğerki hiçbir şekilde sıkıştırılamadıysa dosya normalde karakterin boyutu 1 bytetir fakat Tokenlerin boyutu karakterin boyutunun 3 katına sahip olduğundan dosya boyutu 3 katına çıkmaktadır. Lz77 algoritması kısa metinlerde veya eşleşme olmayan metinlerde işe yaramamaktadır.

B. Kullanılan Bazı Fonksiyonlar

- struct token *Lz77encode:
Programın ana fonksiyonlarından birisidir. Metin.txt içerisinde okunan veriyi parameter olarak alıp token yapıları ile sıkıştırma işlemini gerçekleştirmektedir. Sıkıştırma işlemi gerçekleştirildikten sonra oluşturulan tokenleri return etmektedir.
- char offsetLengthOlustur:
Lz77 yapısında arama tamponu ve ileri tampon olarak iki yapıdan oluşmaktadır. Arama tamponu ve ileri tampondan gelen verileri token içerisinde yer alan 16 bitlik değişkenin içine 12 biti arama tamponu 4 biti ileri tampon olacak şekilde yerleştirmektedir.
- char* fileRead:
Bu fonksiyon kullanıcının Metin.txt içerisine girdiği verileri okuyarak Lz77 nin üzerinde işlem yapması için kaydetmektedir.

C. Karşılaşılan Bazı Sorunlar

- Bu projede karşılaşılan sorunlardan biri token yapısını oluşturmak, arama tamponu ve ileri tampondan gelen bilgileri 16 bit boyutundaki değişkenin içerisine yerleştirmek oldu.
- Lz77 algoritmasının çalışma mantığını anlamakta bazı sorunlar yaşadık fakat kaynakçada belirtilen yerlerde sorununuzu çözecek bilgiler bulduk.

D. Kazanımlar

- C de pointer yapısının kullanımını daha iyi anladık.
- C de bit düzeyinde işlem yapmayı ve stdint.h kütüphanesinden gelen değişkenleri ile çalışmayı öğrendik.
- Lz77 algoritmasının nasıl bir yöntem ile sıkıştırma yaptığını öğrendik.

E. Çalışma Adımları

bin	8.04.2020 01:50	Dosya klasörü	
obj	8.04.2020 01:50	Dosya klasörü	
main.c	16.05.2020 14:43	C Dosyası	6 KB
main.exe	16.05.2020 14:43	Uygulama	33 KB
main.o	16.05.2020 14:43	O Dosyası	4 KB
Metin.txt	16.05.2020 14:28	Metin Belgesi	1 KB
prolab5.cbp	8.04.2020 01:25	CBP Dosyası	2 KB
prolab5.depend	14.05.2020 17:48	DEPEND Dosyası	1 KB
prolab5.layout	15.05.2020 01:49	LAYOUT Dosyası	1 KB

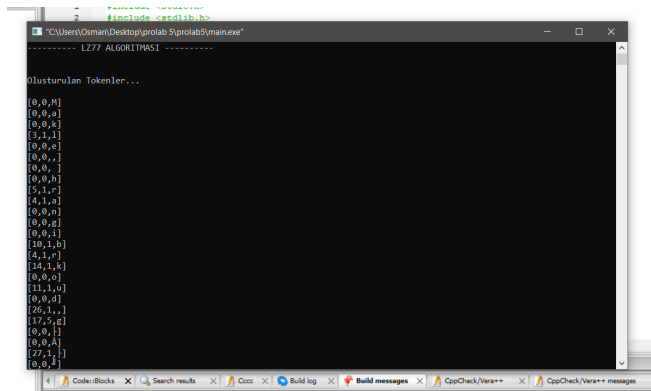
bin	8.04.2020 01:50	Dosya klasörü	
obj	8.04.2020 01:50	Dosya klasörü	
main.c	16.05.2020 14:43	C Dosyası	6 KB
main.exe	16.05.2020 14:43	Uygulama	33 KB
main.o	16.05.2020 14:43	O Dosyası	4 KB
Metin.txt	16.05.2020 14:28	Metin Belgesi	1 KB
prolab5.cbp	8.04.2020 01:25	CBP Dosyası	2 KB
prolab5.depend	14.05.2020 17:48	DEPEND Dosyası	1 KB
prolab5.layout	15.05.2020 01:49	LAYOUT Dosyası	1 KB

Adım 1:

Proje veri sıkıştırma projesi olduğundan dolayı bir metin dosyasına girilen yazıların sıkıştırılması beklenmektedir. Bu yüzden sistem dosyalarının olduğu yerdeki 'Metin.txt' dosyasının içerisine sıkıştırmak istediğiniz metni girmeniz gerekmektedir.

Adım 2:

Metin.txt dosyasının içerisine verileri girdikten sonra program çalıştırılmak istendiğinde main.c dosyasına girerek 'F9' tuşuna basarak codeblock aracılığı ile çalıştırabilirsiniz ya da main.c dosyasının hemen altında yer alan main.exe ye tıklayarak programın direk çalışması sağlanabilir.

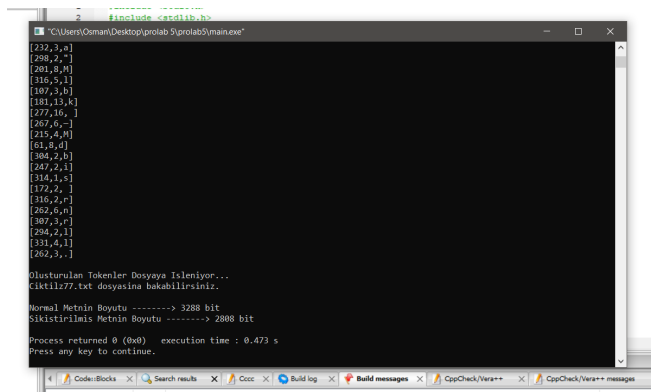


Adım 3:

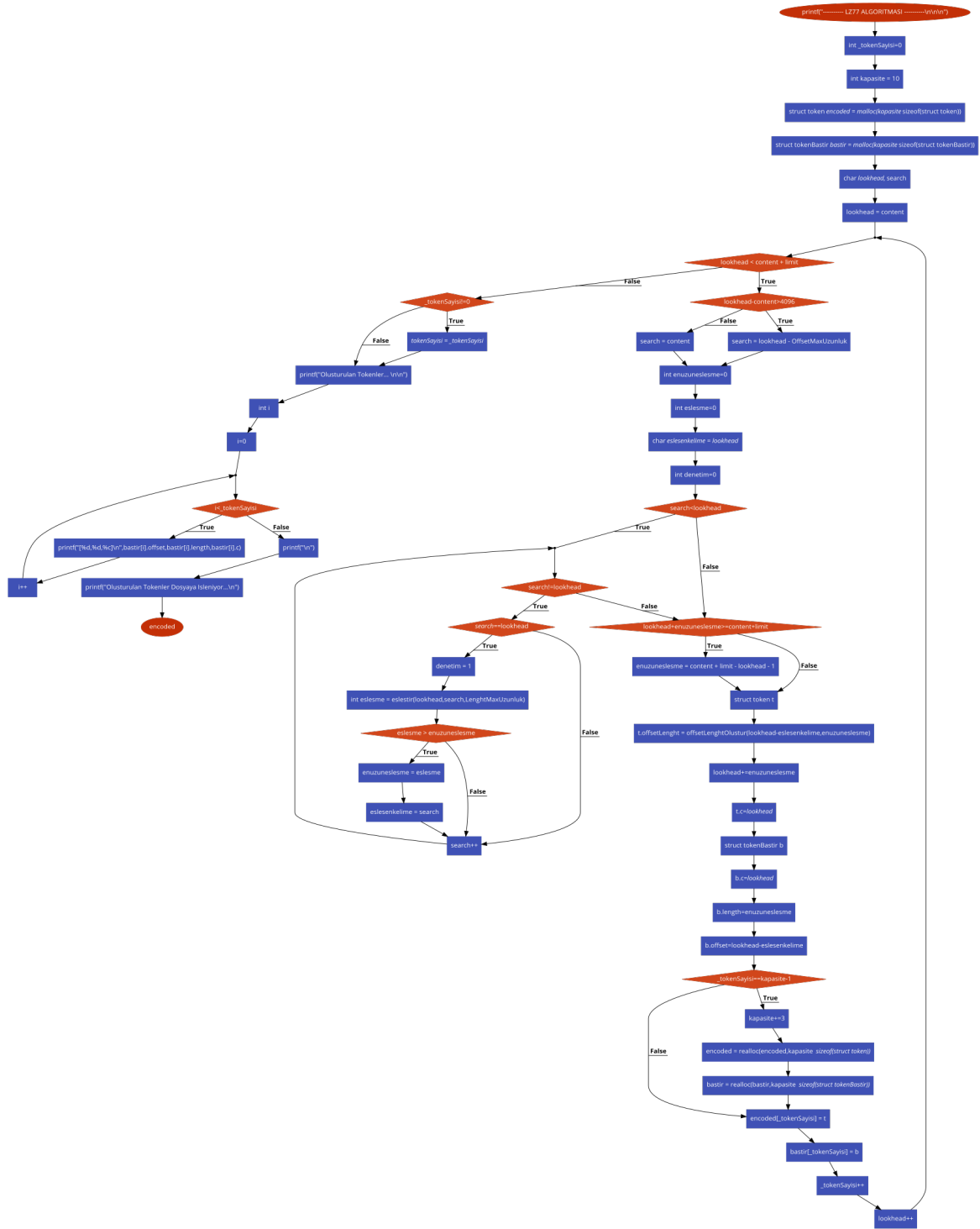
Yandaki resimde de görüldüğü üzere program çalıştırıldığında karşımıza cmd ekranı gelmektedir. Bu ekranda dosyanın sıkıştırıldığı bilgisi verilmektedir. Lz77 ye göre sıkıştırılan verilerden oluşan tokenler ekrana basılmaktadır.

Adım 4:

Dosya sıkıştırılması tamamlandığı vakit cmdde ekranının en aşağıya doğru indiğinizde sıkıştırma sonucu oluşan tüm tokenler gösterilmekte. En aşağıda sıkıştırmanın bittiği bilgisi ekranda gösterilmektedir. Sıkıştırma tamamlandığında normal metnin boyutu ile sıkıştırılan metnin boyutu ekrana basılmaktadır. Sıkıştırma tamamlandıktan sonra otomatik olarak proje dosyalarının yer aldığı yere 'Ciktiliz77.txt' dosyası oluşturularak sıkıştırılan metin içerisine basılmaktadır. Projeye sonlanınca cmd ekranında sonda kullanıcıdan bir girdi beklemektedir. Bu cmd ekranının kapanmaması için alınan bir önlemdir. Projeyi sonlandırdıktan sonra 'Ciktiliz77.txt' dosyasını kontrol edebilirsiniz.



F. Akış Diyagramı



KAYNAKÇA

- [1] http://altanmesut.trakya.edu.tr/pubs/dr_tez.pdf
- [2] https://tr.qwe.wiki/wiki/LZ77_and_LZ78
- [3] <https://www.youtube.com/watch?v=FvWQsuwbraA>
- [4] <http://ysar.net/algoritma/lz77.html>
- [5] <https://www.slideshare.net/veysiertekin/lz77-lempelziv-algorithm>