

Benzer ve Aynı Dosyalar

“Herkes tek ve biriktir, diğerk herkes gibi” bu deęilse de buna oldukça benzer bir laftı ve bu yazı için gzel bir giriş cmlesiydi, ben de gerekeni yaptım. Bilgisayarımızda aynı dosyaların “isteęimiz dıřında” bulunmasından (edebi versiyonu için “at kořturmasından”) rahatsızlık duyarız (duyuyorsunuz deęil mi). Bu dosyalar yazılar, resimler, sesler ve videolar olabilir.

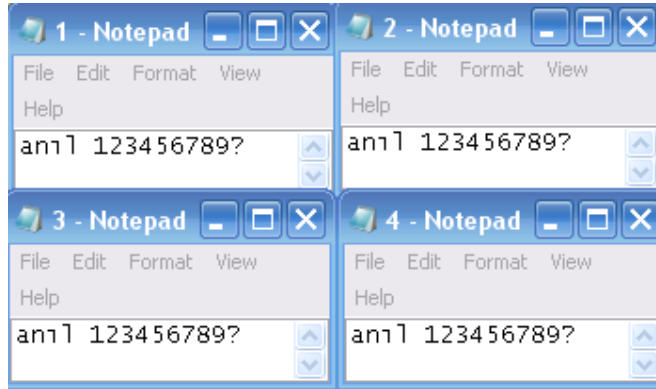
Aynı dosyaların canımızı sıkma derecesi bizim için anlamlı olup olmadıęı ile yakından ilgilidir. Örneęin “C:\Program Files\Common Files\Microsoft Shared\OFFICE12\Office Setup Controller\OSETUP.DLL” dosyasının “C:\MSOCache\All Users\{90120000-0030-0000-0000-0000000FF1CE}-C\osetup.dll” veya herhangi başka bir yerdeki “osetup.dll” ile aynı olması bizi pek bir ilgilendirmez. Anlamlı dosyalarımız içinde ise aynı dosyalara asla katlanamayız, arřivcilięe aykırıdır çünkü. Aynı dosyaları bulmak için yardımcı olan birçok program olmasına raęmen, işimize yarayacak “benzer dosyaları bulucu” programlar pek yoktur. Kısıtlı sayıdaki programlar ise yalnızca tek bir alana yoęunlařmaktadır.

Aynılık ve benzerlik kavramlarını örnekler üzerinde inceleyelim ilk önce:



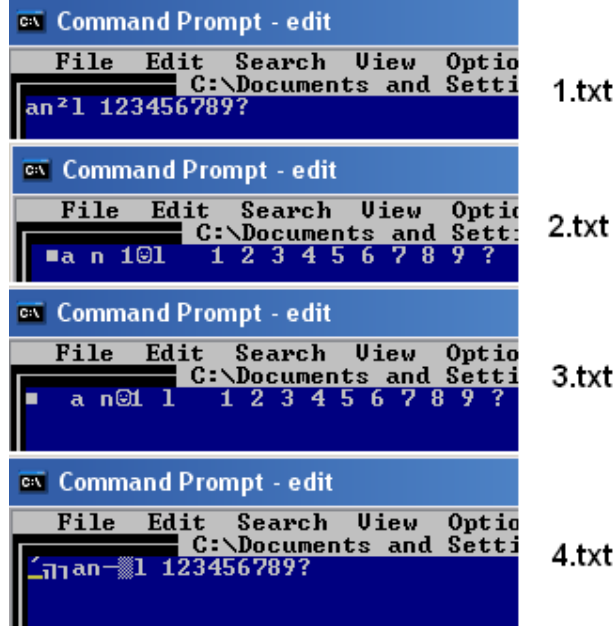
Resim 1. 1 kb lık dört yazı dosyası

Resim 1 deki dosyalar aynı olabilir, çünkü boyutları aynı. Ama bir dosyayı sadece dıř görünüşüyle tanıyamayız, içindekileri de görerek onu gerçekten tanırız.



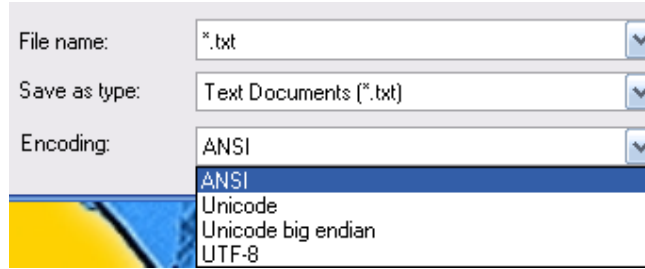
Resim 2. Aynı içerikli dört yazı dosyası

Resim 2 deki dosyalar için artık “aynı olabilir” şphencilięinden, “bu dosyalar aynıdır” kesinlięine geçebilir miyiz? Çoęunlukla dikkatli bakmayız ve yanlış fikirlerle gereksiz zaman harcarız. Bilgisayarda hiçbir “ama” kalmayana kadar emin olabiliriz, denemeye devam edelim.



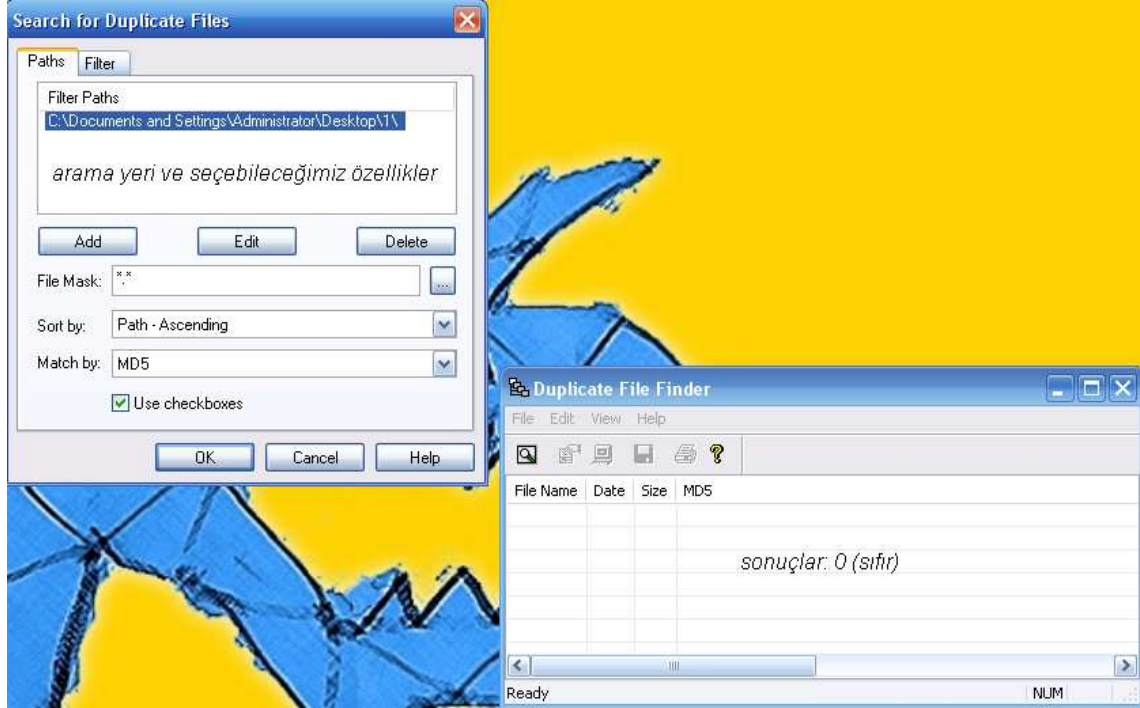
Resim 3. Daha içeriden bir kesit

Dışı bir olsa bile içi farklı gördüğümüz gibi. Bir yazı dosyasını bile farklı kod standartlarıyla kaydettiğimizde yalnızca kendine benzer olduğu için, aynı veya benzer içerikli dosyaların bulunması normal programlama teknikleriyle (veya programlarla) oldukça zordur.



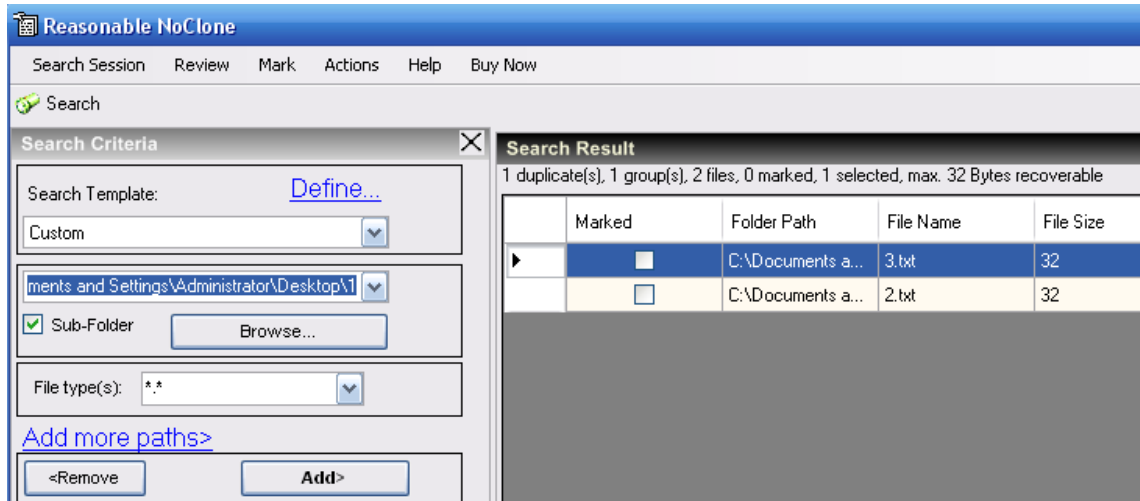
Resim 4. İşin sihri

Örneğin “Dubicate File Finder” isimli program bu dört tane dosyanın aynı olmadığı üzerine yemin ediyor, ya ona inanacaksınız ve mutlu bir şekilde yaşamınıza devam edeceksiniz ya da beni takibe devam edin.



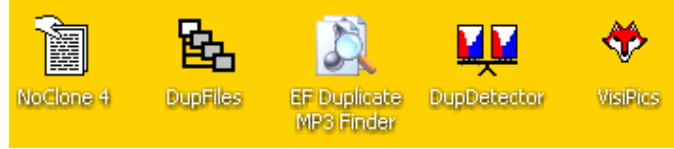
Resim 5. Birebir aynı dosyaları bulan program

Aynı dosyalar üzerinde "hain" deneylerimi şiddetle sürdürdüm. Bu sefer "Reasonable NoClone" isimli programı kullandım. Kullanıcıya daha fazla seçenek sunduğu için (yalnızca binary karşılaştırma yerine, aynı boyutlu dosyaları da bulma özeliği gibi) "2.txt" ve "3.txt" nin aynı olduğunu bulabildi.



Resim 6. Boyutları aynı dosyaları da bulabilen program

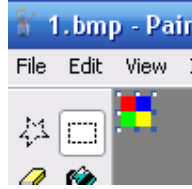
İlk önce kısa bir mazeret yazısıyla giriş yapayım bu paragrafa; internet erişimim kısıtlı imkanlar dâhilinde olduğundan inceleme yapabildiğim programlar ve veriler küçük bir kümeyi kapsar durumda kaldı ama yazının ileriki sürümlerinde (yalana bak!) araştırmayı daha da genişletebilirim. Çünkü eminim ki bir yerlerde bir pdf ile doc un içeriğini karşılaştırabileceğimiz programlar mevcut, aramaya inanmamız kâfi. İncelediğim programlar:



[Reasonable NoClone](#) [Dubicate File Finder](#) [EF Duplicate MP3 Finder](#) [Dup Detector](#) [VisiPics](#)

Resim 7. Ağıma düşürdüğüm programlar

İçeriğe göre arama hakkında hazırladığım lafları sunmadan önce birkaç resim daha göstermek istiyorum sizlere konunun daha iyi anlaşılmasına katkısı olacağını düşünerek.



Resim 8. 2x2 lik bir resim

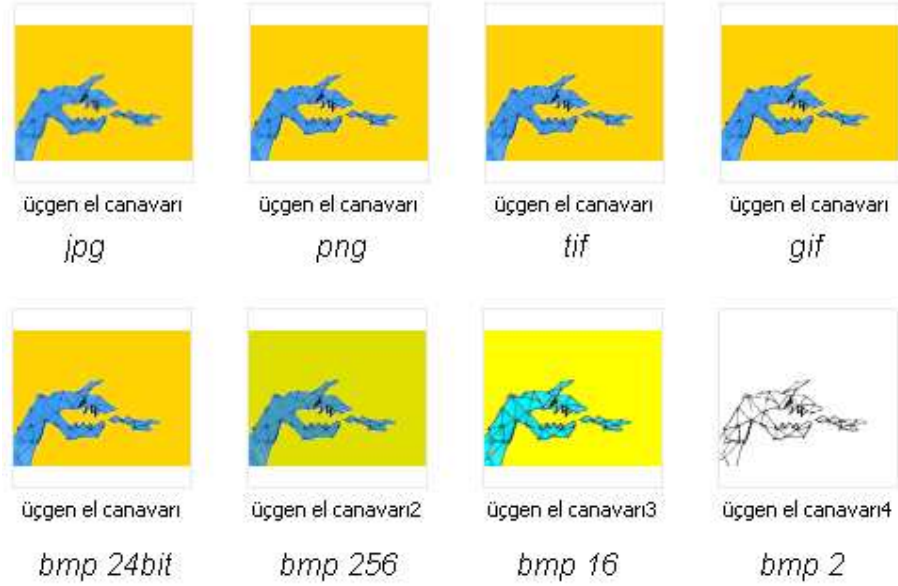
1.bmp nin içeriği:

```
BMF_____6____(_____Ä_____Ä_____ÿ_ÿÿ_____ÿÿ_____
```

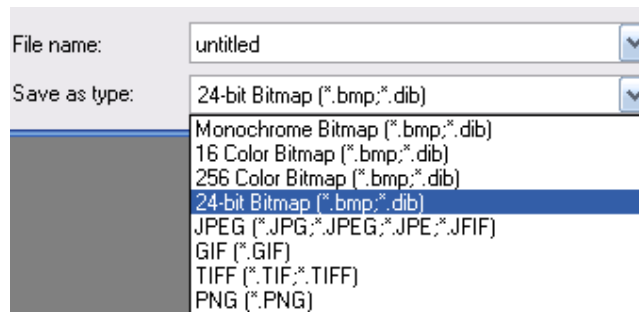
1.gif in içeriği:

```
GIF89a____÷_____€_____€€_____€€_____€€€€€ÄÄÿ_____ÿ_ÿÿ_____ÿÿ_ÿ_ÿÿÿÿÿ_____
_____3_f_____î_____
ÿ_3____33_3f_3™_3î_3ÿ_f____f3_ff_f™_fî_fÿ_____™3_™f_™™™_™î_™ÿ_î_____î3_îf_î™_îî_îÿ_ÿ
__ÿ3_ÿf_ÿ™_ÿî_ÿÿ3____3_33_f3_™3_î3_ÿ33_33333f33™33î33ÿ3f_3f33ff3f™3fî3fÿ3™_3™
33™f3™™3™™î3™ÿ3î_3î33îf3î™3îî3îÿ3ÿ_3ÿ3ÿf3ÿ™3ÿî3ÿÿf____f_3f_ff_™f_îf_ÿf3_f33f3
ff3™f3îf3ÿff_ff3fffff™ffîffÿf™_f™3f™ff™™f™îf™ÿfî_fî3fîffî™fîîîÿÿÿ_ÿ3fÿÿffÿ
™fÿîfÿÿ™_____3™_f™_™™™_î™_ÿ™3_™33™3f™3™™3î™3ÿ™f_™f3™ff™f™™fî™fÿ™™™_™™™3™™™f™™™™™™™™
î™™™ÿ™î_™î3™îf™î™™™îî™îÿ™ÿ_™ÿ3™ÿf™ÿ™™™ÿî™ÿÿî_____î_3î_fî_____îî_ÿî3_î33î3fî3™î3îî3
ÿîf_îf3îffîf™îfîîfÿî™_î™3î™fî™™î™îî™ÿîî_îî3îîfîî™îîîîÿîÿ_îÿ3îÿfîÿ™îÿîîÿÿÿ_
_ÿ_3ÿ_fÿ_™ÿ_îÿ_ÿÿ3_ÿ3ÿ3fÿ3™ÿ3îÿ3ÿÿf_ÿf3ÿffÿf™ÿÿfîÿÿÿ™_ÿ™3ÿ™fÿ™™ÿ™îÿ™ÿÿî_ÿî
3ÿîfÿî™ÿîîÿîÿÿÿ_ÿÿ3ÿÿfÿÿ™ÿÿîÿÿÿ!ù_____,_____ _¹µ0¢/ _;
```


Örneğin bir resmin farklı çözünürlüğe sahip kopyalarını iga (içeriği göre arama) yöntemleriyle kolayca bulabiliriz. (Dup Detector ve VisiPics %100 başarı gösterdi bu testte) (Programcı arkadaşlar burada “nasıllar hakkında” biraz açıklama yaparlarsa biz “meraklı tüketicileri” sevindirirler. Ben kendi tahminlerimi yazayım, dosyalar hızlı bir formülle ortak bir çözünürlüğe çevrilir ve karşılaştırma yapılır veya resimlerin örneğin yukarıdan %20 içeri, soldan %25 gibi bölgelerindeki renk kümelerine bakarak karşılaştırma yapılabilir,...) Bir insan için “çocuk oyuncağı” olan bu işi yapabilen yalnızca iki programın olması bence ne anlam ifade ediyor kısmı sonuç kısmında yer alacak.

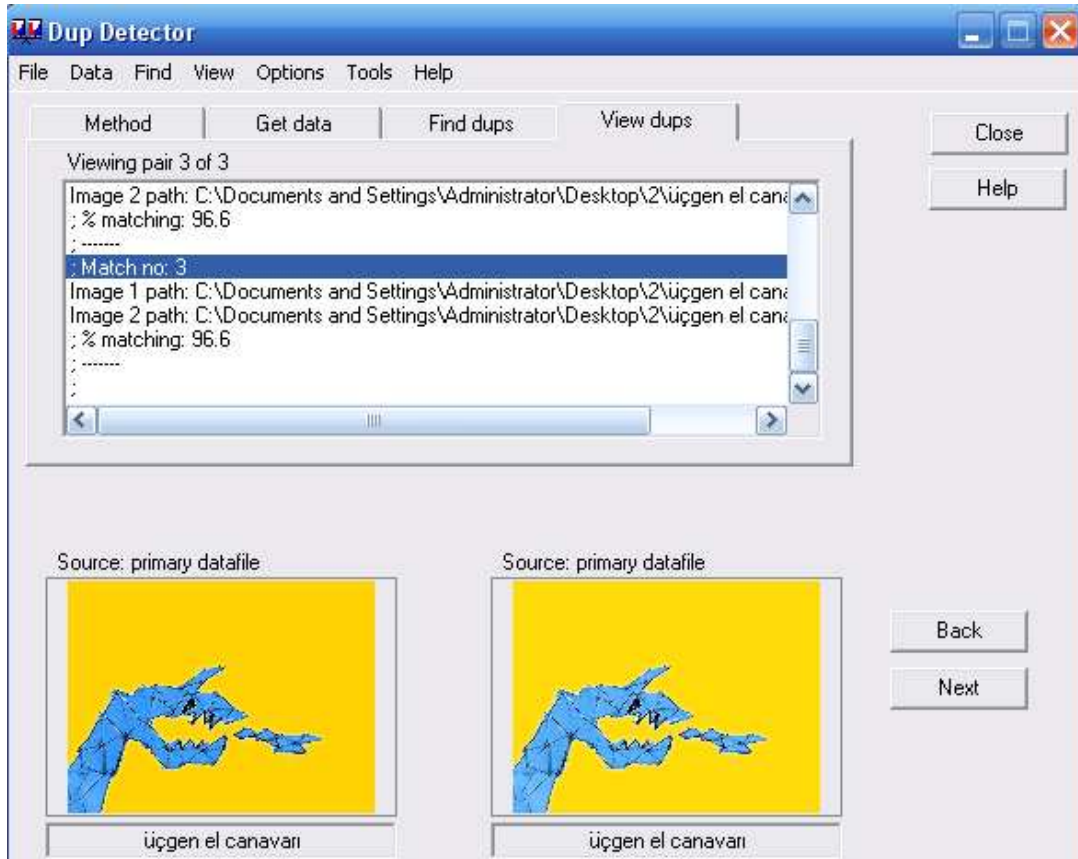
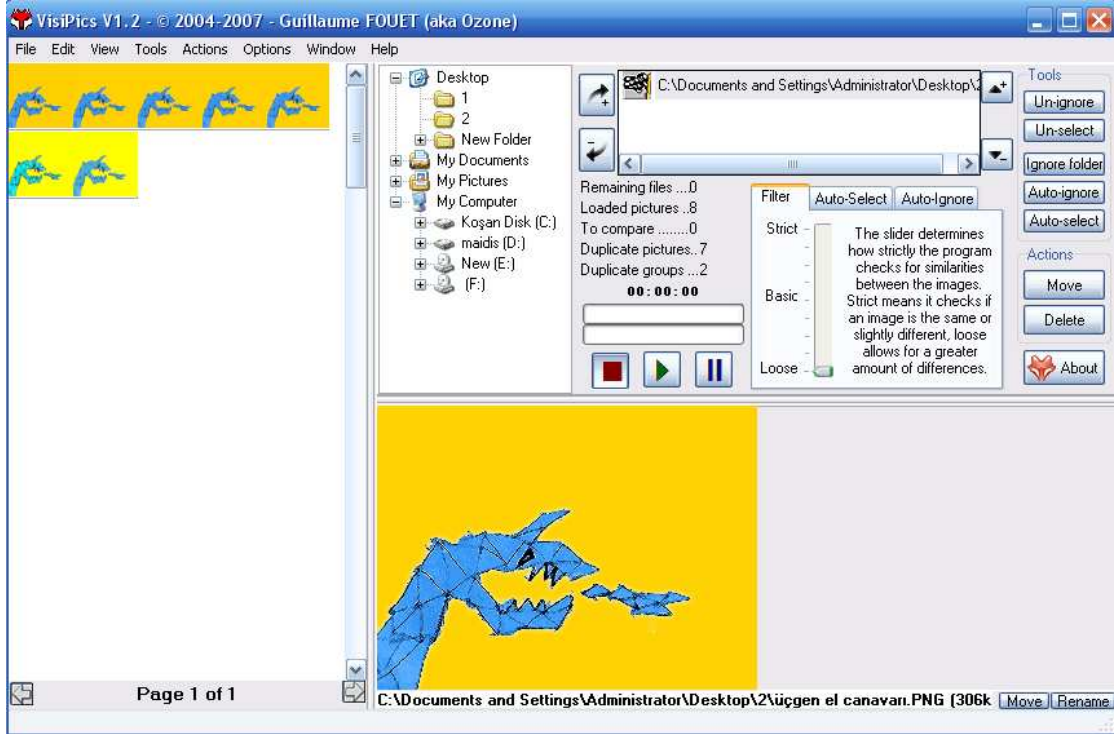


Resim 10. Aynı dosya farklı formatlarda



Resim 11. Bazı görüntü formatları

Sonraki testte bir dosyayı farklı formatlarda kaydettim ve yaptığım test sonuçları aşağıdaki gibidir. (Dup Detector bu sefer tatmin edici sonuçlar veremedi, VisiPics ise görevini başarıyla tamamlamanın mutluluğunu duyumsayamadıysa da benden bir aferin aldı)



Resim 12. Test sonuçları



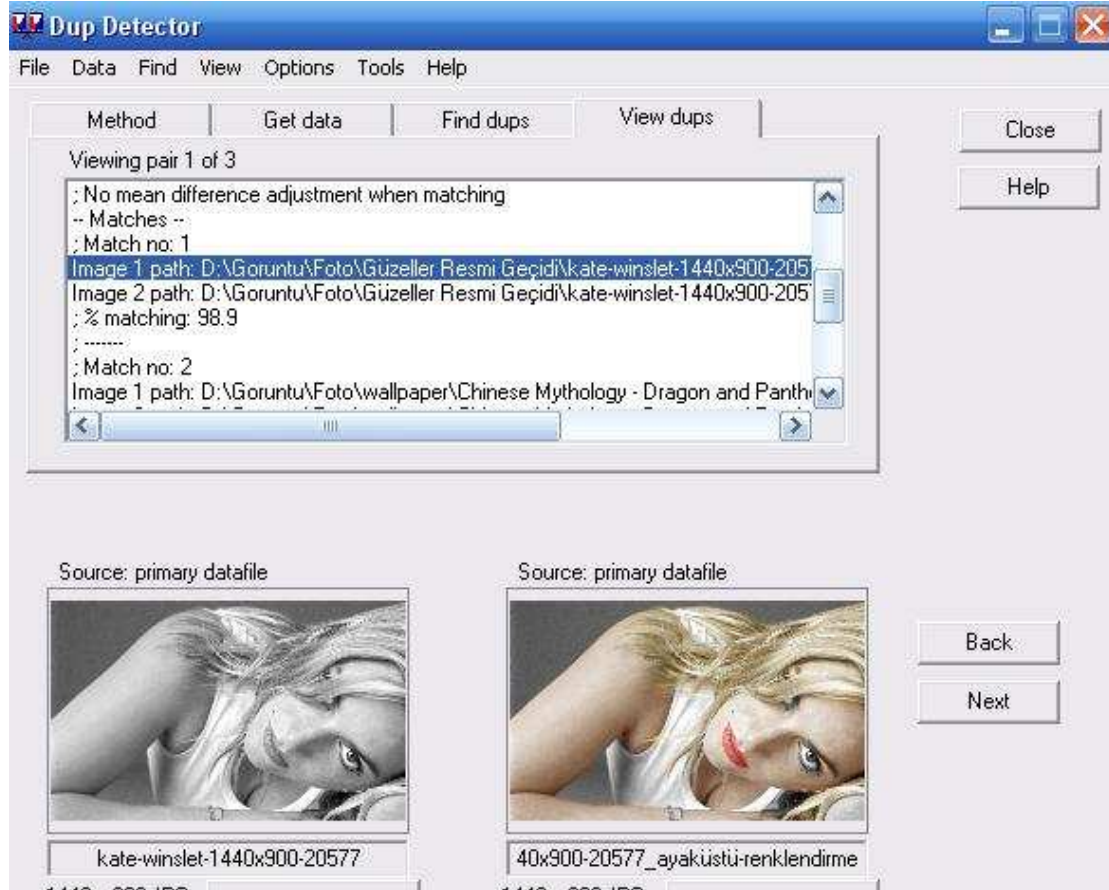
*fotoğrafın orjinal ismi
mb-fp-bw-1.jpg dir sahibini
bilmiyorum*

*sonradan gelen edit: sanırım
run rola run*

Resim 13. Bazı bubi tuzakları

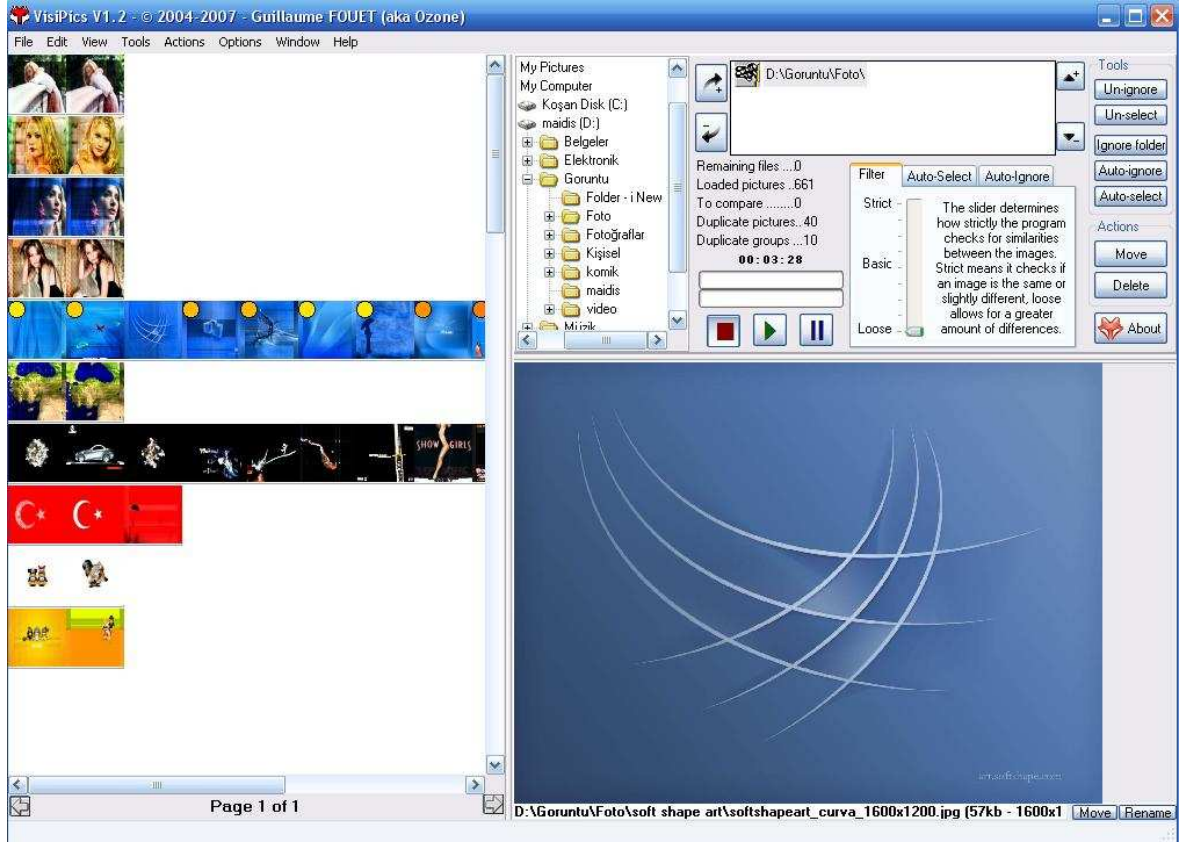
Daha da sonra bir resmin yönünü değiştirerek, aynalayarak ve renklerinin tersini alarak çeşitli denemeler yaptım, Dup Detector bu sefer sınıfı zar zor geçerken, VisiPics hiçbir varlık gösteremedi.

Bir sonraki testte ise internetten indirdiğim ve içeriğini genel olarak tanımlayabileceğim “foto” klasöründe arama gerçekleştirdim. Sonuçlar:



Resim 14. Dup Detector ün bulduklarından

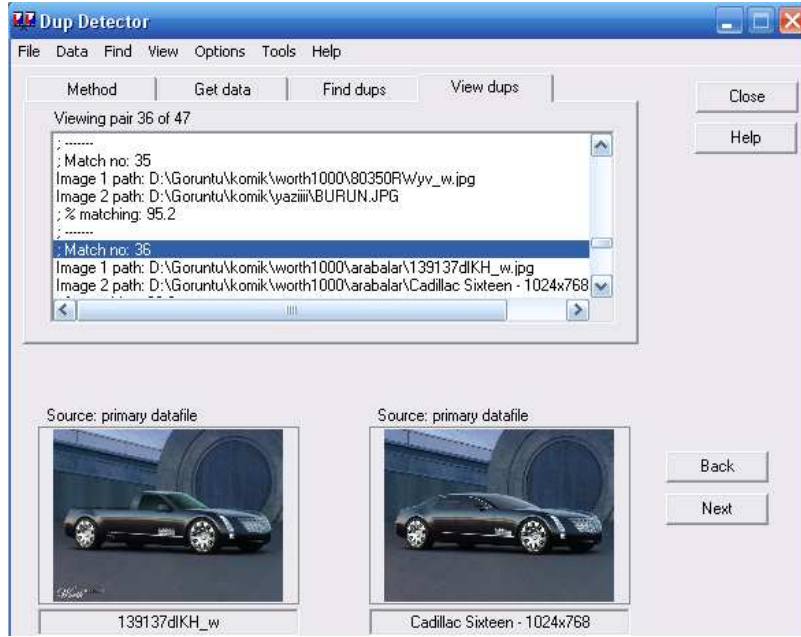
Aynı fotoğrafın renkli ve siyah beyaz kopyasını başarıyla bulan Dup Detector, genelde moda çekimlerinde karşımıza çıkan seri fotoğraflarda oldukça silik kaldı



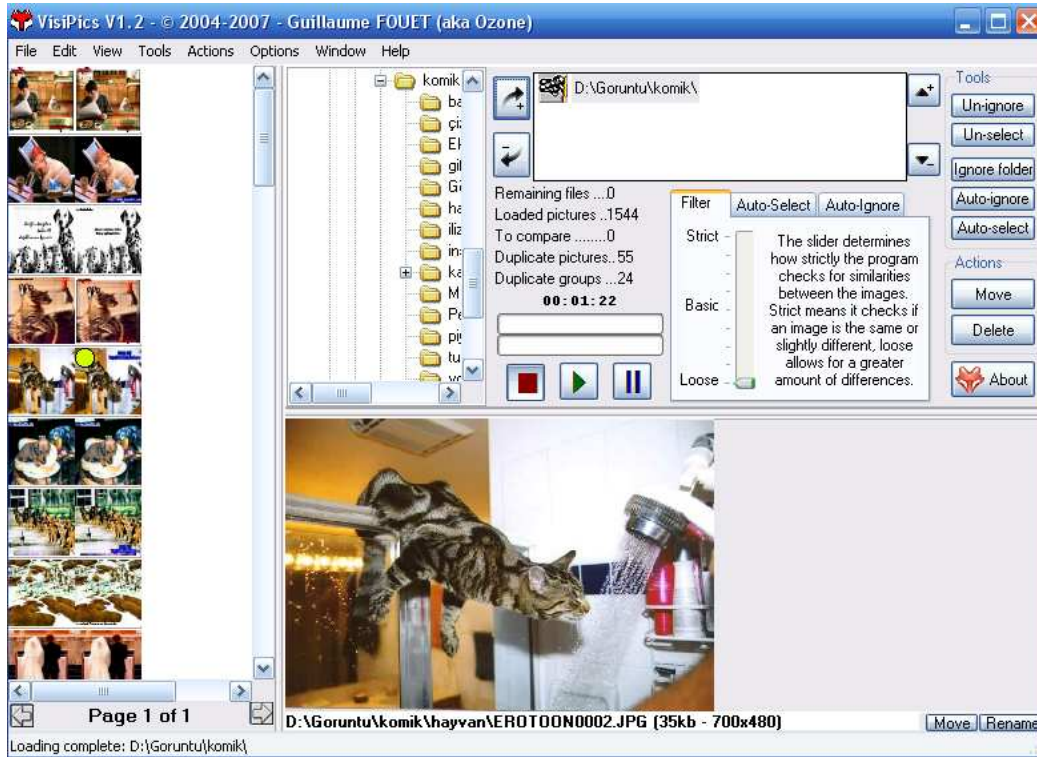
Resim 15. VisiPics in bulduklarından

Bu klasörde harikalar yaratan VisiPics, bulması gereken tüm dosyaları listeleterek gönüllerdeki yerini biraz daha sağlamlaştırdı.

Birçok bilgisayar kullanıcısının bilgisayarının bir yerlerinde konumlandığı ve içinde karikatür, montaj fotoğrafları gibi eserlerin yer aldığı "komik" ya da benzer bir isimde klasörü vardır. Şimdiki testimiz tam da burada.

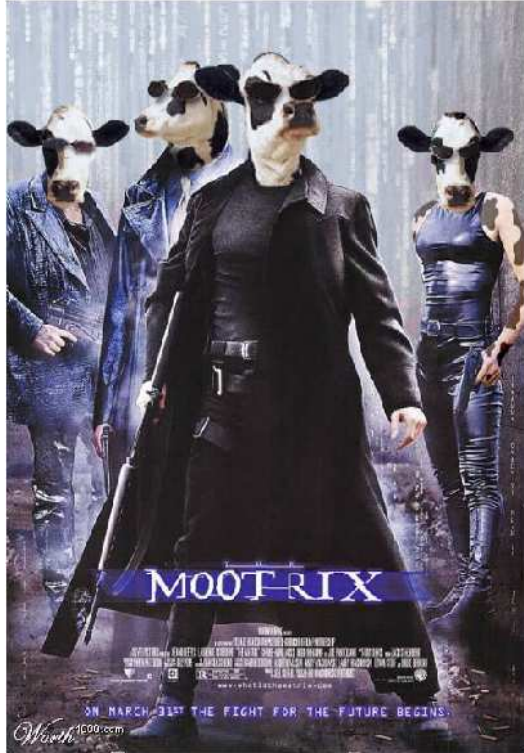


Resim 14. Hanım koş Dup Detector bir şeyler buldu

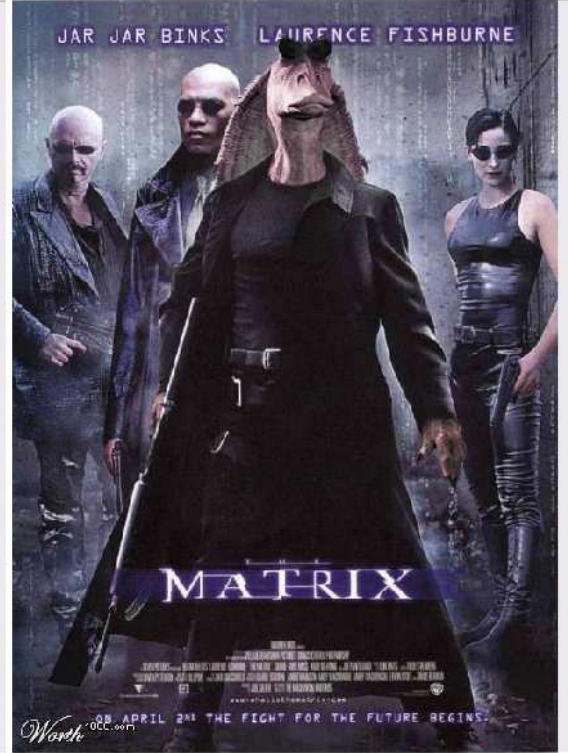


Resim 14. Kedidir kedi

Bu testi açık ara önde olarak VisiPics kazandı. Hiçbir dosyayı isklamadı. VisiPics in bulduklarına daha yakından bir bakış için:



D:\Goruntu\komik\worth1000\film afiş\130806RwAV_w.jpg (66kb - 500x722)



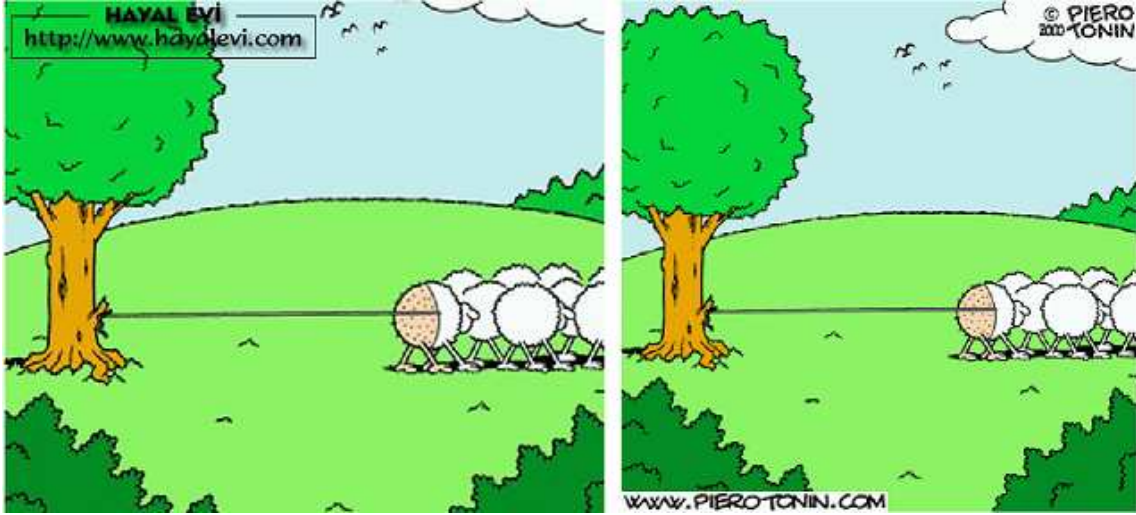
D:\Goruntu\komik\worth1000\film afiş\105095cHKG_w.jpg (56kb - 500x691)



D:\Goruntu\komik\worth1000\insan\69706HmOL_w.jpg (46kb - 500x646)



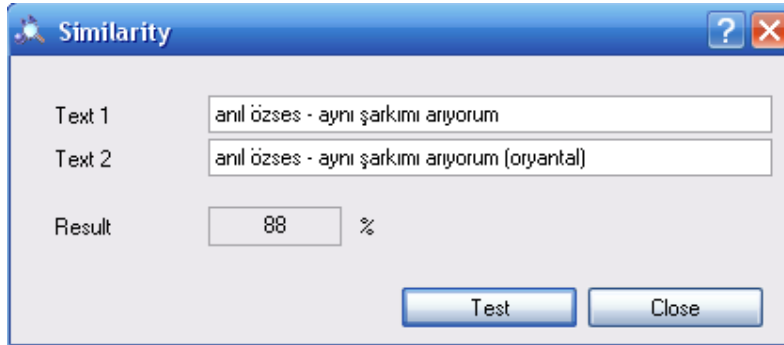
D:\Goruntu\komik\worth1000\insan\69699nsUR_w.jpg (44kb - 500x644)



Resim 15. Tıpkısının ayısının biraz farklısı

Tek karelik görüntüler için tatminkâr aramalar yapılabilmektedir. Bu teknolojilerin biraz daha geliştirilmesiyle hareketli görüntüler için de kolayca uyarlanabilir (Örneğin hareketli bir görüntünün kendi içinde yeterli farklılığa sahip (tamamen aynı renk değil, vs...) bir karenin alınarak diğerleriyle karşılaştırılması zaman alıcı gibi görünse de ilk akla gelen yöntem oldu).

Ses dosyaları için ise durum bu kadar parlak değil ya da ben yanlış yerlerde aradım hep. Demek istediğim ses analizi yaparak benzer dosyaları bulan program yok (yeteri kadar araştırma yapmadığım için bu lafımı yemeye seve seve razıyım). Klasik birebir arama yapan programlar haricinde şarkıcı ismi, parça ismi gibi anahtar ifadelerle bakarak arama yapan programlar mevcut sadece.



Resim 16. Hangi yöntem

Aynı ses dosyalarını bulan arama programlarındaki “şimdilik” en işe yarar yöntem şarkıcı ismi ve parça ismi benzer olan (bkz: resim 16) dosyaları eşleştirmektir. Aynı dosyaları bulması ise sadece arşivinizdeki mp3 bilgilerinin doğru girilmiş olmasına ve şansa bağlıdır.

Sonuç kısımları her zaman zorlandığım bölümler olmuştur. Kısa keseceğim. Çoğu kişinin şu anda ihtiyaç bile duymadığı bu yazılımlar ileride çok önemli bir konumda olacaklardır. Bilgi beraberinde bilgi kirliliğiyle birlikte hızla artmaktadır. Bir makaleyi veya şiiri okuyarak birbirinden esinlenip esinlenilmediğinin tespiti, fotoğraf ve resimleri inceleyerek taklit edilip edilmediği gibi şeyler ve daha birçoğunun çıkış noktası olacağını düşünüyorum bu konunun.

Maidis