

# Üç Boyutlu El Hareketlerinin Tanınması için İki Boyutlu Özniteliklerin Birleştirilmesi

## Combination Strategies for 2D Features to Recognize 3D Gestures

Oya Aran, Lale Akarun

Bilgisayar Mühendisliği Bölümü  
Boğaziçi Üniversitesi, İstanbul

aranoya@boun.edu.tr, akarun@boun.edu.tr

### Özetçe

Bu çalışmada, iki kameradan oluşan bir kurulumda üç boyutlu el hareketlerini tanıyan bir sistem tasarladık. Kameraların kalibrasyon matrislerinin hesaplanmadığı ve üç boyutta geri çatma işleminin yapılamadığı durumlarda, iki boyutlu el gezinmesinin her iki kamera için ayrı ayrı çıkartılarak öznitelik ya da karar aşamasında birleştirilmesi başarımın büyük ölçüde artmasını sağlamaktadır. İki boyutlu el gezinmeleri, elin orta noktası ve eli çevreleyen elipsin uzunluk, genişlik ve açısı izlenerek çıkarılmış ve bu gezinmeler Kalman Süzgeci kullanılarak yumuşatılmıştır. Boyutları ve başlangıç noktaları düzgelene gezinmeler Saklı Markov Modelleriyle (HMM) veya sabit uzunlukta yeniden örneklenerek Destek Vektör Makinalarıyla (SVM) sınıflandırılmışlardır. Farklı kameralardan çıkarılan gezinmelerin farklı aşamalarda birleştirilmesinin performans üzerindeki etkisi incelenmiştir. En iyi sonuç, farklı kameralardan çıkarılmış gezinmelerin HMM kullanılarak modellenmesi ve karar aşamasında birleştirilmesiyle elde edilmiş, bu yöntemle toplam 210 test örneğinde sadece 1 ya da 2 örnekte hata yapılmıştır.

### Abstract

In this study, using a two camera setup, we designed a system that recognizes 3D gestures. When 3D reconstruction is not possible or infeasible, combining 2D hand trajectories at feature or decision level increases the system performance drastically. The trajectories are extracted by tracking the center-of-mass of the hand and the width, height and orientation of the enclosing ellipse. Trajectories are then smoothed using a Kalman Filter. Following the translation and scale normalization, the trajectories are modelled using Hidden Markov Models (HMM) and using Support Vector Machines (SVM) by converting the trajectories to fixed length using re-sampling. Trajectories extracted from different cameras are combined at different levels and the effect to the system performance is observed. The best result is obtained by modelling the trajectories using HMMs and combining at decision level, with %1 error in 210 test examples.

### 1. Giriş

El hareketleri insanlar arası iletişimde kimi zaman sadece konuşmayı tamamlamak için kimi zaman ise tek başına bir

iletişim aracı olarak kullanılmaktadır. Son dönemde, bilgisayar ve kamera teknolojisindeki yenilikler el hareketlerinin aynı zamanda insan bilgisayar iletişiminde de diğer araçlarla beraber ya da tek başına kullanılabilmesini göstermiştir [1]. El hareketlerinin tanınması aynı zamanda işaret dilinin tanınması açısından da önemlidir. El hareketi, şekli ve konumunun analizine dayanan işaret dilinde yüz ifadesi, baş-vücut hareketleri gibi bileşenler el bilgisini destekleyerek işaret dilini zenginleştirirler [2].

El hareketi tanıma sistemleri elin uzay-zamansal bileşenlerini modellerler. Uzaysal bileşenin karmaşıklığı uygulamaya göre değişiklik gösterir. En karmaşık halinde elin ve parmakların yapısı detaylı şekilde incelenmelidir. Daha basit olarak genel el şekli veya sadece elin pozisyonu kullanılabilir. Zamansal bileşen ise iki ya da üç boyutta el gezinmesinin çıkarılmasıyla elde edilir. Gezingenin çıkarılmasında çeşitli el izleme teknikleri kullanılabilir gibi (Kalman Süzgeci [3], Parçacık Süzgeci [4], vs.) zamansal şablonlara [5] dayalı teknikler de kullanılabilir. Çıkarılan gezinme Sonlu Durumlu Makineler, Zaman Gecikmeli Sinir Ağları, Saklı Markov Modelleri ya da şablon eşleştirme teknikleriyle modellenebilir [6]. Bu teknikler içinde, Saklı Markov Modelleri, bir çok sistemde başarı sağlamışlardır.

Saklı Markov Modelleri, aynı sınıfa ait aynı ya da farklı uzunluktaki örnekleri kullanarak o sınıf için bir model oluştururlar. Her sınıf için farklı bir Saklı Markov Modeli oluşturulur ve bir test örneği geldiğinde olabilirliği en yüksek olan model / sınıf seçilir. Farklı uzunluklardaki örnekleri kullanma ve sıralı veriyi modellemedeki başarısı Saklı Markov Modellerini el hareketi tanıma problemi için uygun bir teknik haline getirmektedir. Fakat sınıflandırma problemlerinde, sınıf içi benzerlikleri olduğu kadar sınıflar arasındaki ayrımları da gözönüne alan tekniklerin daha başarılı olduğu bilinmektedir. Bu tür teknikler arasında Destek Vektör Makinaları, Yapay Sinir Ağları sayılabilir.

Bu çalışmada, farklı kameralardan çıkarılan iki boyutlu el gezinmeleri ve genel el şekli bilgisi öznitelik aşamasında ya da karar aşamasında birleştirildi ve iki farklı sınıflandırıcı (HMM, SVM) kullanılarak sınıflandırıldı. İki boyutlu el gezinmelerini öznitelik aşamasında ya da karar aşamasında birleştirmenin üç boyutlu el hareketlerinin tanınmasına katkısı incelendi.

## 2. İki Boyutlu El Gezinesinin Çıkarılması

Bu çalışmada yedi el hareketinden oluşan bir veri tabanı kullanıldı [7]. Veri tabanındaki el hareketleri iki el birden kullanılarak yapılan ve üç boyutlu nesnelere itirmek ve döndürmek için kullanılacak hareketleri içermektedir. Kullanıcılar her iki ellerine farklı renklerde (sol elde mavi ve sağ elde sarı) eldivenler giymişlerdir. Kullanıcının sağ ve soluna iki kamera yerleştirilmiştir. Eğitim kümesinde, her kişiden 10 adet olmak üzere, dört kişiden alınmış toplam 280 örnek bulunmaktadır. Test kümesinde ise farklı üç kişiden alınmış toplam 210 örnek bulunmaktadır.

El gezinesinin çıkarılması için aşağıdaki adımlar sırasıyla uygulanır:

1. Elin bulunması (Giyilen eldivenin rengine göre eşikleme ve bağlantılı bileşen algoritması kullanılarak)
2. Bulunan iki boyutlu koordinatların Kalman Süzgeci kullanılarak yumuşatılması
3. Ölçek farklılıklarının düzelenmesi ( $x, y$  koordinatlarının iki boyutlu düzelenmesiyle)
4. Konum farklılıklarının düzelenmesi (Gezinenin başlangıç noktasını (0,0) koordinatına çekerek)

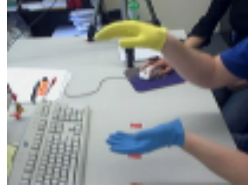
### 2.1. Elin Bulunması ve İzlenmesi

Kullanıcının her iki elinin bulunması için ilk olarak giyilen eldivenlerin renklerine göre eşikleme uygulandı (Şekil 1b). Eşiklenen imgeler bağlantılı bileşen algoritmasıyla bölümlendi. Bölütleme sonucunda oluşan imgede yer alan en fazla alana sahip olan bileşenin el olduğu kabul edildi (Şekil 1c).

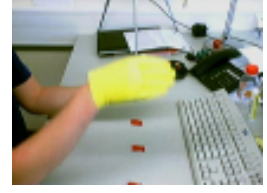
El gezinesi, her çerçevede bulunan elin ağırlık merkezi kullanılarak çıkarılmış ve çıkarılan gezinje Kalman Süzgeci kullanılarak yumuşatılmıştır. Gezinede, ağırlık merkezine ek olarak eli çevreleyen elipsin uzunluğu, genişliği ve açısı öznitelik olarak kullanılmıştır. Bu sayede genel el şekil bilgisi el hareket bilgisiyle birleştirilmiştir.

Bu kurulumda kameraların kullanıcıya yandan bakıyor olması sebebiyle, eller bazı hareket sınıflarında birbirlerini kapatabilirler. Bu kapatma ancak kameraların birinde olabilir. Aynı anda iki kamerada birden ellerin birbirini kapatması söz konusu değildir. Şekil 2 sağ kamerada sağ elin sol eli kapattığı bir anı göstermektedir. Aynı anda sol kamera görüntüsünde ise herhangi bir kapatma yoktur. El bulma algoritmasının varsayımı, eşiklenmiş imgede bulunan en fazla alana sahip bağlantılı bileşenin el olduğu varsayımdır. Bu varsayım doğru olduğu sürece, eller birbirini kısmen kapatsa da el bulma işlemi, Kalman Süzgeci tarafından yumuşatılabilen ufak bir hata ile gerçekleştirilebilir. Bu varsayım doğru değilse el olduğu varsayılan bileşenin el ile hiç bir ilgisi olmayabilir. Bu bilgi Kalman Süzgecinin parametrelerini güncellemek için kullanılırsa, Kalman Süzgecinin güvenilirliği azalabilir. Bu problemi çözmek için bulunan bileşenin alanı bir eşik değeri ile karşılaştırılır. Eğer alan bu değer altındaysa bir elin diğerini kapattığı varsayılır ve Kalman Süzgecini parametreleri güncellenmez. Kapatmanın gerçekleştiği çerçeve sayısı bir-iki çerçeve ile sınırlı kaldığı durumda, Kalman Süzgecinin yaptığı tahminlerin kullanılması problem yaratmaz.

Sol Kamera



Sağ Kamera



Şekil 2: Kameraların birinde eller birbirini kapatabilir ama aynı anda diğer kamerada iki el birden görülür.

### 2.2. Düzgeleme

Uygulamaya göre değişebilmekle birlikte, aynı el hareketi farklı kişiler tarafından, hatta aynı kişi tarafından, farklı ölçeklerde ve farklı yerlerde yapılabilir. Bu durumun sınıflandırma başarısına etki etmemesi için öteleme ve ölçek farklılıklarının düzelenmesi gerekir.

$((x_1, y_1), \dots, (x_t, y_t), \dots, (x_N, y_N))$  iki boyutlu el gezinesi ve  $N$  gezinje uzunluğu olsun. Öteleme düzgelemesinde kullanmak için,  $x_m$  ve  $y_m$  şu şekilde tanımlanır:

$$x_m = \frac{x_{max} + x_{min}}{2}, \quad y_m = \frac{y_{max} + y_{min}}{2} \quad (1)$$

ve gezinenin  $x$  ve  $y$  koordinatlarındaki orta noktasını gösterirler. Ölçek düzgelemesinde kullanmak için,  $\delta_x$  ve  $\delta_y$  şu şekilde tanımlanır:

$$\delta_x = \frac{x_{max} - x_{min}}{2}, \quad \delta_y = \frac{y_{max} - y_{min}}{2} \quad (2)$$

ve gezinenin  $x$  ve  $y$  koordinatlarındaki saçılımını gösterirler. Gezinenin genel şeklini bozmamak için, ölçekleme faktörü  $x$  ve  $y$  koordinatlarındaki saçılımın maksimumu olarak seçilir:

$$\delta = \max(\delta_x, \delta_y) \quad (3)$$

Bu bilgiler kullanılarak, düzelenmiş gezinje,  $((x'_1, y'_1), \dots, (x'_t, y'_t), \dots, (x'_N, y'_N))$ ,  $0 \leq x'_t, y'_t \leq 1$  olmak üzere, şu şekilde hesaplanır:

$$x'_t = 0.5 + 0.5 \frac{x_t - x_m}{\delta} \quad (4)$$

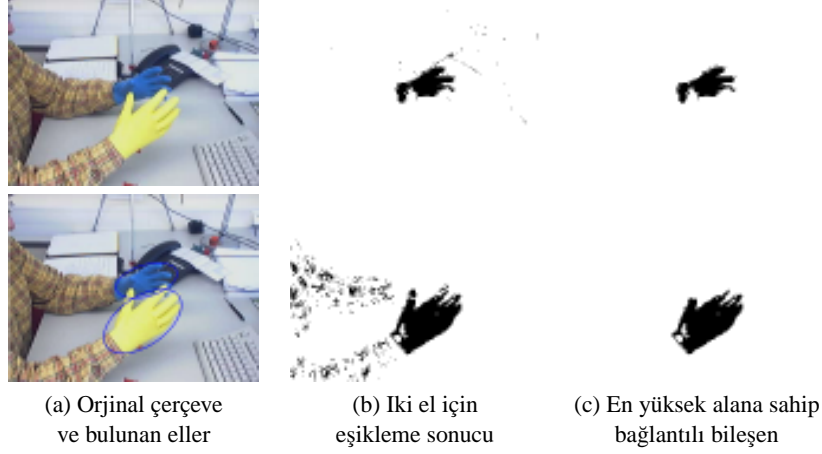
$$y'_t = 0.5 + 0.5 \frac{y_t - y_m}{\delta} \quad (5)$$

Gezinedeki el şekline ait diğer bilgiler de benzer şekilde düzelenmelidir. Elipsin genişlik ve uzunluğu aynı  $x$  ve  $y$  koordinatlarındaki düzgelemede yapıldığı gibi saçılımın maksimumu kullanılarak yapılmalıdır. Aksi halde şekil bilgisi bozulabilir. Elips açısı ise bağımsız olarak düzelenmelidir.

En son olarak, düzelenmiş gezinje başlangıç noktası (0,0) koordinatına çekilir.

## 3. Deneyler

Her el hareketi için dört değişik gezinje çıkarılabilir: birinci kameradan sol ve sağ el gezinje (L1 ve R1) ve ikinci kameradan sol ve sağ el gezinje (L2 ve R2). Ellerin birbirini kapatma ihtimali olduğundan tek bir kameradan



Şekil 1: Çerçeve de elin bulunması

elde edilen gezinme hataları içerebilir. Ayrıca, tek bir kamera bilgisini kullanmak sınıflandırıcıyı iki boyutlu hareket bilgisine sınırlayacaktır. Kalibrasyon matrisinin olmaması sebebiyle üç boyuta geri çatma işlemi yapılamasa da iki kameradan gelen bilgi birleştirilerek üç boyutlu hareket bilgisi bir ölçüde elde edilebilir. İki boyutlu gezinmelerin değişik aşamalarda birleştirilmesinin sistem başarımına etkisini görmek amacıyla değişik deneyler yaptık:

1. L1R1: Birinci kameradan sol ve sağ el gezinmeleri (Çerçeve başına öznitelik vektörü boyutu: 10)
2. L2R2: İkinci kameradan sol ve sağ el gezinmeleri (Çerçeve başına öznitelik vektörü boyutu: 10)
3. L1L2R1R2: Birinci ve ikinci kameradan sol ve sağ el gezinmelerinin öznitelik aşamasında birleştirilmesi (Çerçeve başına öznitelik vektörü boyutu: 20)
4. L1R1  $\oplus$  L2R2: Birinci ve ikinci kameradan sol ve sağ el gezinmelerinin karar aşamasında birleştirilmesi (Çerçeve başına öznitelik vektörü boyutu: 10+10)

Her deney için iki sınıflandırıcı eğitildi: HMM ve SVM. Değişken uzunluktaki gezinmeleri SVM kullanılarak eğitebilmek için yeniden örnekleme yoluyla sabit uzunluğa çevirdik. Gezingenin  $x$  ve  $y$  noktaları, uzamsal yeniden örnekleme yoluyla ve doğrusal aradeğerleme kullanarak 12 noktaya örneklendiler. Şekil 3'de değişik nokta sayılarıyla yeniden örnekleme sonuçları görülebilir. Örneklenmiş noktaların, genel el şeklini yansıtan öznitelikleri ise kendilerine zamansal olarak en yakın noktadan kopyalandı. SVM sonuçları LibSVM kullanılarak alındı [8]. Eğitilen SVM'lerde çekirdek fonksiyonu olarak Radyal Taban Fonksiyonları kullanıldı. Çok sınıflı SVMler bire-bir yöntemiyle eğitildi. Eğitimin öncesinde tüm öznitelikler z-normalizasyon kullanılarak düzelendi ve ortalamadan iki standart sapma dışında kalan değerler iki standart sapmaya çekilerek aykırı değerlerin etkisi azaltıldı. HMMlerde, durum atlama sol-sağ modeli kullanıldı. El hareketlerinin uzunluğuna bakılarak (en uzun 40 çerçeve) dört durumlu bir model ve her durumda bir Gauss bileşeni kullanıldı. Durum sayısını ya da bileşen sayısını arttırmanın başarıma iyi yönde farkedilir bir etkisi olmadığı görüldü.

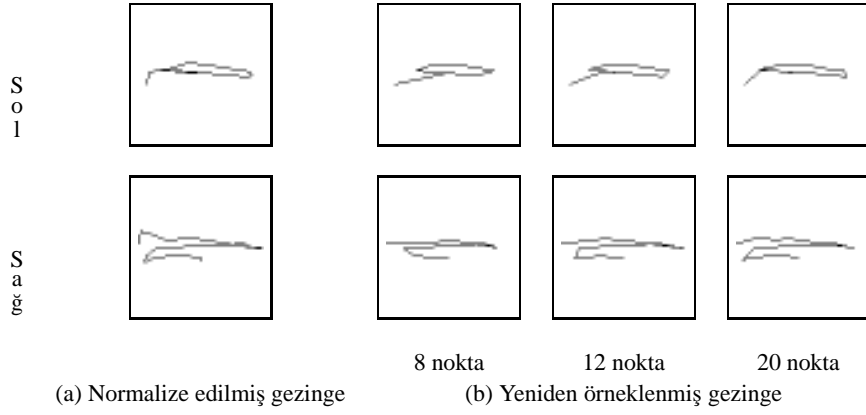
Tek kamera bilgisi kullanılarak elde edilen başarımlar, öznitelik ve karar aşamasında birleştirme sonuçları Tablo 1'de görülebilir. Karar aşamasında birleştirme için, HMMlerin olabilirlik sonuçları ve SVMlerin ardıl olasılıkları kullanıldı ve toplam kuralıyla birleştirildi. Tek kamera bilgisi kullanılarak elde edilen %5 civarındaki hata, iki kameradan alınan bilginin öznitelik aşamasında birleştirilerek beraber kullanılmasıyla %50'den fazla azalarak %2 civarına düşürüldü. Karar aşamasında HMM kullanarak birleştirildiğinde ise hata daha da azalarak %1'e düştü. SVM'ler kullanılarak yapılan birleştirme ise öznitelik aşamasındaki birleştirmeden daha farklı bir sonuç vermedi.

SVM sınıflandırıcısı, L1R1 kümesi hariç, HMM sınıflandırıcısından daha fazla hata yaptı. HMM'lerin sıralı bilgiyi daha iyi modellediği ve yeniden örnekleme sırasında önemli bazı bilgilerin kaybolmuş olabileceği gözönüne alınırsa bu beklenen bir sonuçtu. HMM'den daha fazla hata yapmış olmasına karşın, SVM sonuçlarının HMM sonuçlarına yakın çıkmış olması, değişik uzunluktaki gezinmelerin sabit uzunluğa çevrilmesinde daha iyi bir yöntem kullanıldığı takdirde hatanın azaltılabileceğini gösteriyor.

Tablo 1: Test hataları ve standart sapmalar

Dataset	SVM	HMM
L1R1 ( <i>cam1</i> )	0.057 $\pm$ 0.000	0.060 $\pm$ 0.011
L2R2 ( <i>cam2</i> )	0.062 $\pm$ 0.000	0.054 $\pm$ 0.012
L1R1L2R2	0.024 $\pm$ 0.000	0.022 $\pm$ 0.007
L1R1 $\oplus$ L2R2	0.024 $\pm$ 0.000	<b>0.009 <math>\pm</math> 0.004</b>

Karar aşamasında birleştirme uygulandığında hatanın azalması, birleştirilen sınıflandırıcıların birbirinin hatalarını yüksek oranda düzeltmesi sayesinde gerçekleşir. Kullandığımız veri kümesinde, farklı kameralardan çıkarılan gezinmelerle eğitilen HMM'ler birleştirildiğinde hata, %1'e kadar azaltılabilmektedir. Tablo 2 ve 3'de birinci ve ikinci kamera bilgisine eğitilen HMM'lerin hata matrisleri görülebilir. Birinci kameradan çıkarılan gezinmelerde daha çok ikinci ve beşinci sınıflara ait örnekler hatalı sınıflandırılmış. İkinci kamerada ise ik-



Şekil 3: Değişik nokta sayılarıyla yeniden örnekleme

inci, dördüncü ve yedinci sınıfların örneklerinde hata yapıldığı görülmektedir. Birleştirme sonucunda, sınıflandırıcı sadece ikinci ve dördüncü sınıflarda birer hata yapmaktadır.

Tablo 2: L1R1 HMM hata matrisi

30	0	0	0	0	0	0
0	27	0	3	0	0	0
0	0	29	0	1	0	0
0	0	1	29	0	0	0
0	0	2	0	25	3	0
0	0	0	0	0	30	0
0	0	1	0	0	0	29

Tablo 3: L2R2 HMM hata matrisi

30	0	0	0	0	0	0
0	27	0	2	1	0	0
0	0	29	0	0	0	1
1	0	0	28	0	1	0
0	0	0	0	30	0	0
0	0	0	0	1	29	0
0	0	0	1	0	1	28

#### 4. Sonuçlar

Bu çalışmada üç boyutta yapılan el hareketlerini iki boyut için çıkarılmış özneliklerle tanıyan bir sistem oluşturduk. Üç boyutta geri çatma işlemi yapılmadan, iki kameradan çıkarılan iki boyutlu el gezinmelerinin öznelik ya da karar aşamasında birleştirilmesinin sistemin başarımını arttırdığı görüldü. Sistem, gezinmeler öznelik aşamasında birleştirildiğinde, toplam 210 test örneği için %2.5 hata, karar aşamasında birleştirildiğinde ise %1 hata yapmaktadır.

Bu çalışma DPT/03K120250 "Algısal İnsan Bilgisayar Etkileşimi" ve AB 6. çerçeve programı SIMILAR projeleri tarafından desteklenmektedir.

Tablo 4: L1R1 $\oplus$ L2R2 HMM hata matrisi

30	0	0	0	0	0	0
0	29	0	1	0	0	0
0	0	30	0	0	0	0
0	0	1	29	0	0	0
0	0	0	0	30	0	0
0	0	0	0	0	30	0
0	0	0	0	0	0	30

#### 5. Kaynakça

- [1] Vladimir Pavlovic, Rajeev Sharma, and Thomas S. Huang, "Visual interpretation of hand gestures for human-computer interaction: A review.," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 19, no. 7, pp. 677–695, 1997.
- [2] Sylvie C. W. Ong and Surendra Ranganath, "Automatic sign language analysis: A survey and the future beyond lexical meaning.," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 6, pp. 873–891, 2005.
- [3] R. E. Kalman, "A new approach to linear filtering and prediction problems.," *Transactions of the ASME—Journal of Basic Engineering*, vol. 82, pp. 35–45, 1960.
- [4] M. Isard and A. Blake, "Condensation – conditional density propagation for visual tracking.," *International Journal of Computer Vision*, vol. 26, no. 1, pp. 5–28, 1998.
- [5] A. Bobick and J. Davis, "Real-time recognition of activity using temporal templates.," in *Proceedings of the Workshop on Applications of Computer Vision*, 1996.
- [6] Ying Wu and Thomas S. Huang, "Hand modeling, analysis, and recognition for vision based human computer interaction.," *IEEE Signal Processing Magazine*, vol. 21, pp. 51–60, 2001.
- [7] Sebastian Marcel and Agnes Just, "IDIAP Two handed gesture dataset.," Available at <http://www.idiap.ch/~marcel/>.
- [8] Chih-Chung Chang and Chih-Jen Lin, *LIBSVM: a library for support vector machines*, 2001, Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.