



دانشگاه صنعتی امیرکبیر
دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد
گرایش هوش مصنوعی

تشخیص بی درنگ حرکات دست در تعامل انسان با کامپیوتر

نگارش

سید ناصر نوراشرف‌الدین

استاد راهنما

دکتر رضا صفابخش

اسفند ۸۶

بسمه تعالی



دانشگاه صنعتی امیرکبیر

(بلی تکنیک تهران)

معاونت پژوهشی

فرم اطلاعات پایان نامه

کارشناسی ارشد و دکترا

تاریخ:

پیوست:

نام و نام خانوادگی: سیدناصر نور اشرفالدین دانشجوی آزاد بورسیه معادل

شماره دانشجویی: ۸۴۱۳۱۰۳۷ دانشکده: مهندسی کامپیوتر رشته تحصیلی: هوش مصنوعی

نام و نام خانوادگی استاد راهنما: رضا صفا بخش

عنوان پایان نامه به فارسی: تشخیص بی‌درنگ حرکات دست در تعامل انسان با کامپیوتر

عنوان پایان نامه به انگلیسی: Real-time hand gesture recognition in human-computer interaction

نوع پروژه: کارشناسی ارشد دکترا کاربردی بنیادی توسعه ای نظری

تاریخ شروع: ۸۵/۱۱/۲۰ تاریخ خاتمه: ۸۶/۱۲/۲۵ تعداد واحد: ۶

سازمان تأمین کننده اعتبار:

واژه های کلیدی به فارسی: تعامل انسان-کامپیوتر، حرکتهای ایستا، حرکتهای پویا، بینایی ماشین، مدل رنگ پوست.

واژه های کلیدی به انگلیسی: Human-computer interaction, static gestures, dynamic gestures, computer vision, human skin model.

نظرها و پیشنهادها به منظور بهبود فعالیت های پژوهشی دانشگاه:

استاد راهنما:

دانشجو:

امضاء استاد راهنما: تاریخ:

نسخه ۱: معاونت پژوهشی

نسخه ۲: کتابخانه و به انضمام دو جلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز اسناد و مدارک علمی

تقدیم به پدر و مادر عزیزم.

تقدیر و تشکر

بدینوسیله از راهنمایی ها و زحمات استاد عزیز و گرامی دکتر رضا صفابخش که به طرق مختلف در مسیر انجام این پایان نامه من را یاری نمودند، صمیمانه تشکر و قدردانی می‌نمایم. همچنین از تمامی دوستان عزیزی که در تهیه پایگاه داده این پروژه همکاری داشتند نیز صمیمانه تشکر می‌کنم. این پایان نامه بر اساس قرارداد شماره ۵۰۰/۱۲۳۵۱/ت مورخ ۸۶/۸/۲۰ تحت حمایت مالی مرکز تحقیقات مخابرات ایران انجام شده است.

بدینوسیله اینجانب سیدناصر نوراشرف‌الدین تعهد می‌نمایم که مطالب ارائه شده در این پایان‌نامه حاصل کار پژوهشی و تحقیق اینجانب می‌باشد و قبلاً برای احراز مدرک دیگری ارائه نشده است. رجوع به دست‌آوردهای پژوهشی دیگران که در این پایان‌نامه از آنها استفاده شده، مطابق مقررات ارجاع داده شده است. چنانچه در هر شرایطی این موارد بدرستی رعایت نگردد، دانشگاه مجاز به ابطال پایان‌نامه خواهد بود. کلیه حقوق مادی و معنوی این اثر متعلق به دانشکده مهندسی کامپیوتر دانشگاه صنعتی امیرکبیر می‌باشد.

نام و نام‌خانوادگی دانشجو: سیدناصر نوراشرف‌الدین

تاریخ:

امضاء:

چکیده

در سالهای اخیر تحقیقات زیادی برای استفاده از ابزارها و روش های جدید برای ارتباط با کامپیوتر انجام شده است. در حال حاضر برای برقراری ارتباط با وسایلی مانند کامپیوتر، رباتها و یا سایر وسایل الکترونیکی نظیر تلویزیون از ابزارهای واسطی مانند موشواره و صفحه کلید، اهرمهای مکانیکی و الکتریکی و یا کنترلرهای از راه دور استفاده می شود. با پیشرفت فناوری در آینده نزدیک شاهد خواهیم بود که انسانها بسیاری از کارهای روزمره خود را با استفاده از این وسایل انجام می دهند و بنابراین اگر بتوانند با این وسایل به همان شیوه ای که انسانها با یکدیگر تعامل دارند، ارتباط برقرار کنند، کارها بسیار راحت تر و بدون نیاز به ابزارهای واسط قابل انجام خواهد بود.

در این پایان نامه یک الگوریتم مبتنی بر بینایی ماشین ارائه شده است که می تواند بصورت بی درنگ حرکت های دست کاربران را تشخیص داده و معادل دستوری آنها را در کامپیوتر اجرا کند. در ابتدا با استفاده از روش تفاضل پس زمینه، تصویر کاربر از زمینه جدا و سپس با استفاده از یک مدل گاسی، پیکسلهایی که به رنگ پوست انسان هستند مشخص می شوند. از ویژگی های مهم الگوریتم پیشنهادی که برتری آن نسبت به سایر الگوریتم های ارائه شده نیز بشمار می رود، عدم نیاز به داشتن تصویر زمینه از قبل می باشد. در این روش الگوریتم خود از روی فریمهای اولیه تصویر زمینه را می سازد و در فریمهای بعدی آن را بروزرسانی می کند. در مرحله بعد با استفاده از عملگرهای ریخت شناسی نویزهای تصویر و ناحیه های آویزان دست حذف می شود. پس از تقطیع دستها، کف دست از بقیه قسمت های آن مانند ساعد و بازوها جدا می شود. برای شناسایی دستها با استفاده از گشتاورهای مرکزی نرمال شده و گشتاورهای پیرامون، بردار ویژگی هر دست استخراج شده و این بردارهای ویژگی به عنوان ورودی به شبکه های LVQ داده می شوند که این شبکه ها نوع حرکت کاربر را تشخیص می دهند. کلیه حرکت های این سامانه ایستا بوده و کاربران حرکتها را با دو دست انجام می دهند. این الگوریتم می تواند ۲۵ حرکت مختلف را که توسط پنج وضعیت پایه انجام می شود را با سرعت ۲/۵ فریم در ثانیه شناسایی کند. الگوریتم ارائه شده می تواند در حالت برون-خط با دقت ۹۷/۶ و در حالت برخط با دقت ۹۴ درصد حرکت های کاربران را تشخیص دهد. در این الگوریتم مشکل همپوشانی صورت و دستهای کاربر حل شده و برای نور محیط هیچ نورپردازی خاصی در نظر گرفته نشده و سامانه می تواند در نورپردازی های مختلف کار کند.

واژگان کلیدی: تعامل انسان-کامپیوتر، حرکت های ایستا، حرکت های پویا، بینایی ماشین، مدل رنگ پوست.

فهرست مطالب

فصل اول- مقدمه	۱
فصل دوم- بررسی کارهای انجام شده	۱۴
۱-۲ مدل دستها	۱۵
۲-۲ تصویر برداری از حرکت‌های دست	۱۷
۳-۲ تقطیع تصاویر و استخراج ویژگی	۲۳
۴-۲ شناسایی حرکت دستها	۲۸
۵-۲ الگوریتمی برای تشخیص الفبای زبان اشاره ای آمریکایی	۲۸
۶-۲ الگوریتمی برای تشخیص الفبای زبان اشاره ای چینی	۳۹
۷-۲ الگوریتمی تطبیقی برای تشخیص حرکت‌های ایستا	۴۵
۸-۲ الگوریتمی برای تشخیص الفبای زبان اشاره ای عربی	۴۹
۹-۲ الگوریتمی برای تشخیص الفبای زبان اشاره ای باغچه بان	۵۳
۱۰-۲ تقطیع تصویر دست	۵۸
۱۱-۲ استخراج ویژگی	۶۹
۱۲-۲ شناسایی (دسته بندی) دستها	۷۰
۱۳-۲ خلاصه	۷۷
فصل سوم- الگوریتم پیشنهادی	۷۸
۱-۳ دوربین و کاربر	۸۰
۲-۳ تشخیص نواحی به رنگ پوست	۸۱
۳-۳ تقطیع تصاویر دست	۸۵
۴-۳ استخراج ویژگی	۱۰۴
۵-۳ دسته بندی حرکتها	۱۰۶
۶-۳ ارزیابی الگوریتم پیشنهادی	۱۰۷
۷-۳ خلاصه	۱۱۱
فصل چهارم- نتیجه گیری و پیشنهادات	۱۱۲
مراجع	۱۲۰

فهرست شکلها

- شکل ۱-۱ یک نمونه دستکش داده ای ۵
- شکل ۲-۱ استفاده از دستکشهای رنگی به جای دستکشهای داده ای ۶
- شکل ۱-۲ دیاگرام عملیاتی یک سامانه HCI مبتنی بر بینایی ماشین [۱۶] ۱۴
- شکل ۲-۲ مدلهای حرکت دست [۱۷] ۱۵
- شکل ۳-۲ یک مدل اسکلتی دست [۱۷] ۱۶
- شکل ۴-۲ دو حرکت تقریباً مشابه ۱۷
- شکل ۵-۲ چرخش درون صفحه ای دستها ۲۰
- شکل ۶-۲ چرخش خارج از صفحه ای دستها ۲۰
- شکل ۷-۲ اعمال محدودیتها در گرفتن تصویر دستها [۱۴] ۲۳
- شکل ۸-۲ یک نمونه از حرکتهای پویا [۷] ۲۴
- شکل ۹-۲ دو حرکت پویای تقریباً مشابه [۷] ۲۴
- شکل ۱۰-۲ معماری سامانه آقای گوپتا [۲] ۲۹
- شکل ۱۱-۲ دیاگرام عملیاتی سامانه آقای گوپتا [۲] ۳۰
- شکل ۱۲-۲ نمونه هایی از تصاویر گرفته شده برای اعداد صفر تا نه [۲] ۳۰
- شکل ۱۳-۲ نمونه ای از تغییرات درون دسته ای (دسته عدد صفر) [۲] ۳۱
- شکل ۱۴-۲ تقطیع تصاویر با استفاده از روش آتسو [۲] ۳۲
- شکل ۱۵-۲ تصاویر $S_{1,5}$ و $S_{3,5}$ بعد از حذف نویز [۲] ۳۴
- شکل ۱۶-۲ نحوه محاسبه مولفه های LCS [۲] ۳۵
- شکل ۱۷-۲ LCS های شکل‌های $f_{1,5}$ و $f_{3,5}$ [۲] ۳۶
- شکل ۱۸-۲ ۳۰ حرکت زبان اشاره ای چینی [۴] ۳۹
- شکل ۱۹-۲ چند نمونه از تصویرهای استفاده شده در این الگوریتم [۴] ۴۴
- شکل ۲۰-۲ نه حرکت ایستا برای تغییر اسلایدها [۵] ۴۶
- شکل ۲۱-۲ پیکربندی سامانه [۵] ۴۷
- شکل ۲۲-۲ پیدا کردن مکان میج با استفاده تغییر پهنای دست [۵] ۴۸
- شکل ۲۳-۲ ۳۰ الفبای زبان اشاره ای عربی [۳] ۵۱
- شکل ۲۴-۲ تصویر دست به همراه دو مشخصه مرکز ثقل و جهت اصلی [۳] ۵۱

- شکل ۲-۲۵ تصویر دست، پیرامون استخراج شده و پیرامون هموار شده با استفاده از فیلتر گاسی [۳]..... ۵۱
- شکل ۲-۲۶ بردارهای مرکز ثقل تا پیرامون دستها [۳]..... ۵۲
- شکل ۲-۲۷ نحوه محاسبه زاویه بین نقاط پیرامونی..... ۵۵
- شکل ۲-۲۸ دیاگرام عملیاتی تقطیع رنگ پوست [۲۲]..... ۵۹
- شکل ۲-۲۹ معماری شبکه های RCE [۲۶]..... ۶۵
- شکل ۲-۳۰ تقریب توزیع رنگ پوست با استفاده از شبکه RCE [۲۶]..... ۶۶
- شکل ۲-۳۱ تقطیع تصویر دستها با استفاده از روشهای سنجش عمق [۳۵]..... ۶۸
- شکل ۲-۳۲ دستکش داده ای سایبرگلاو [۳۶]..... ۶۹
- شکل ۲-۳۳ تطبیق گراف کشسان با تصویر ورودی [۴۷]..... ۷۳
- شکل ۲-۳۴ نمونه ای از چرخشهای برون صفحه ای [۴۸]..... ۷۴
- شکل ۲-۳۵ دستکشهای رنگی استفاده شده در این سامانه [۴۸]..... ۷۵
- شکل ۲-۳۶ نمونه های آموزشی در زوایای دید متفاوت [۴۱]..... ۷۶
- شکل ۳-۱ پنج وضعیت پایه ۷۸
- شکل ۳-۲ دیاگرام عملیاتی الگوریتم پیشنهادی این پایان نامه ۷۹
- شکل ۳-۳ دو نمونه از تصاویر زمینه ۸۱
- شکل ۳-۴ یکی از کاربران سامانه با لباس آستین کوتاه..... ۸۱
- شکل ۳-۵ نحوه یادگیری وزنه های کوهنن در فضای نمونه های ورودی..... ۸۳
- شکل ۳-۶ هیستوگرام سه بعدی نمونه های رنگ پوست..... ۸۳
- شکل ۳-۷ ماسک گاسی در دو حالت مختلف..... ۸۵
- شکل ۳-۸ یک نمونه ماسک گاسی ۸۶
- شکل ۳-۹ دو حرکت مشابه متوالی ۸۷
- شکل ۳-۱۰ کاهش تغییرات نورپردازی، کاهش نور محیط در تمام صحنه..... ۸۹
- شکل ۳-۱۱ کاهش تغییرات نورپردازی، افزایش نور محیط در تمام صحنه..... ۸۹
- شکل ۳-۱۲ کاهش تغییرات نورپردازی، کاهش نور در بخشی از تصویر..... ۸۹
- شکل ۳-۱۳ کاهش تغییرات نورپردازی، افزایش نور در بخشی از تصویر..... ۹۰
- شکل ۳-۱۴ نتیجه اعمال الگوریتم تفاضل پس زمینه برای چند فریم نمونه..... ۹۱
- شکل ۳-۱۵ نتیجه الگوریتم تفاضل پس زمینه و مدل گاسی برای چند فریم نمونه..... ۹۲
- شکل ۳-۱۶ حذف نویزهای تصویر..... ۹۳

- شکل ۳-۱۷ یک نمونه از همپوشانی دست و صورت کاربر ۹۳
- شکل ۳-۱۸ دو ناحیه متصل در حالتی که یکی از دستها با صورت همپوشانی دارد ۹۴
- شکل ۳-۱۹ فضای جستجو در روشهای تطبیق بلوکی [۵۱] ۹۶
- شکل ۳-۲۰ محاسبه بردار سرعت یکی از بلوکهای فریم جاری [۵۱] ۹۶
- شکل ۳-۲۱ مثالی از الگوریتم تطبیق بلوکی با جستجوی لگاریتمی [۵۱] ۹۸
- شکل ۳-۲۲ تقطیع صورت کاربران از فریمهای ابتدایی ۹۹
- شکل ۳-۲۳ حذف تصویر صورت کاربر برای حذف همپوشانی صورت و دستها ۱۰۰
- شکل ۳-۲۴ حذف نویزهای پیرامونی دستها و پر کردن حفره ها ۱۰۱
- شکل ۳-۲۵ هیستوگرام پهنای چند تصویر دست ۱۰۲
- شکل ۳-۲۶ چند خروجی الگوریتم برش مچ ۱۰۳
- شکل ۴-۱ استخراج فریمها از فیلم ویدئویی حرکت کاربران ۱۱۲
- شکل ۴-۲ تصویرهای خروجی مراحل مختلف الگوریتم (۱) ۱۱۴
- شکل ۴-۳ تصویرهای خروجی مراحل مختلف الگوریتم (۲) ۱۱۵
- شکل ۴-۴ تصویرهای خروجی مراحل مختلف الگوریتم (۳) ۱۱۶
- شکل ۴-۵ تصویرهای خروجی مراحل مختلف الگوریتم (۴) ۱۱۷
- شکل ۴-۶ تصویرهای خروجی مراحل مختلف الگوریتم (۵) ۱۱۸

فهرست جداول

- جدول ۱-۲ نتایج ارزیابی الگوریتم ارائه شده در ۵۰ آزمایش [۲]..... ۳۸
- جدول ۲-۲ بررسی و مقایسه پنج الگوریتم تشخیص حرکتهای دست..... ۵۶
- جدول ۱-۳ مقایسه الگوریتم پیشنهادی با پنج الگوریتم تشخیص حرکتهای دست..... ۱۰۸

فصل اول

۱- مقدمه

از ابتدای تولید کامپیوترها یکی از مهمترین فعالیتهای طراحان این بوده است که کاربران بتوانند هر چه ساده‌تر با کامپیوترها ارتباط برقرار کرده و به راحتی از قابلیت‌های آن استفاده کنند. در ابتدا کاربران تنها از طریق وسائل مکانیکی خاصی با کامپیوترها تعامل داشتند و بتدریج با پیشرفت فناوری ابزارهای جدیدی نظیر موشواره، صفحه کلید و قلم نوری نیز در اختیار کاربران قرار گرفت. متداول‌ترین شیوه تعامل با کامپیوترهای فعلی که برای چندین دهه است که تغییرات زیادی نیز نداشته، استفاده از سخت افزارهایی نظیر موشواره و صفحه کلید می‌باشد. با وجود سادگی این وسایل، هنوز هم افراد زیادی هستند که در استفاده از کامپیوترها مشکل داشته و حتی استفاده از این ابزارها برای کاربران حرفه‌ای نیز در بعضی از موارد خسته‌کننده می‌باشد. امروزه با پیشرفت فناوری شاهد هستیم که کامپیوترها در هر ابعادی به وسایل جانبی نظیر دوربین، میکروفون و سنسورها مجهز هستند و از طرف دیگر آمار و ارقام نیز نشان می‌دهد که در آینده‌ای نزدیک انسانها بیشتر وقت خود را با کامپیوترها می‌گذرانند تا با سایر انسانها. بنابراین این نیاز احساس خواهد شد که انسانها بتوانند طوری با کامپیوترها ارتباط برقرار کنند که گویی با انسانهای دیگر تعامل دارند.

از طرف دیگر با توجه به افزایش روز افزون نقش رباتها در زندگی مدرن امروزی و گسترش زمینه‌های کاربرد آنها به منظور تسهیل شرایط زندگی، ایجاد زمینه و راهکارهای مناسب برای برقراری ارتباط با آنها به بهترین و راحت‌ترین روش ممکن از جمله چالش‌های مهم در فناوریهای امروزی است. نیازهای بشر امروزی با توجه به افزایش روزافزون فناوری، انتظارات بیشتری را در این شیوه ارتباطی ایجاد کرده است. بنابراین یک فناوری پیشرفته در طراحی یک ربات شدیداً تحت تاثیر ارتباطی است که محصول مورد نظر، به عنوان یک وسیله هوشمند با کاربر خود برقرار خواهد کرد. هر چقدر شیوه برقراری این ارتباط طبیعی‌تر بوده و به نحوه ارتباط مابین انسان‌ها نزدیک‌تر باشد، ارتباط مورد نظر مناسب‌تر و راحت‌تر ارزیابی می‌شود و کاربران تمایل بیشتری برای استفاده از این محصول خواهند داشت. برای تحقق این موضوع لازم است که کامپیوترها (یا رباتها) مشابه انسانها عمل کرده و کارهای پایه‌ای نظیر دیدن، تشخیص صدا، صحبت کردن و لمس کردن را یاد بگیرند. همچنین همانطور که ما انسانها در برخورد با رفتارهای متفاوت به شیوه‌های مختلفی عمل می‌کنیم، لازم است که کامپیوترها نیز بتوانند کاربران متفاوت خود را مدل کرده و با هر کاربر مطابق نیازهای آن رفتار کنند.

امکان بینایی برای کامپیوترها باعث می‌شود که این دستگاهها بتوانند جهان اطراف خود را ببینند و تغییرات آن را درک کنند. این قابلیت کامپیوترها را قادر خواهد کرد که حرکتهای کاربران خود را حس کنند و پس از تجزیه و تحلیل تصاویر، حرکتهای کاربران را فهمیده و دستورها و پیغامهای موجود در تصاویر را استخراج و اجرا کنند. برای تولید اینگونه واسطه‌های^۱ مبتنی بر بینایی، آشنایی با مباحث بینایی ماشین، پردازش تصویر، گرافیک کامپیوتری، علوم زبان شناسی و روانشناسی مورد نیاز است. با وجود اینکه پیشرفتهای زیادی در این مباحث و علوم صورت گرفته، اما تعامل مبتنی بر بینایی^۲ پیشرفتهای زیادی نداشته است و مشکلات زیادی برای تولید یک واسطه کاربری مناسب وجود دارد. یکی از مهمترین مشکلاتی که برای تولید این واسطهها وجود دارد این است که حرکت انسانها دارای پیچیدگی‌های زیادی است و الگوریتم جامعی برای مدل کردن تمامی حرکتها وجود ندارد. همچنین میزان تغییرات حرکتها در انسانهای مختلف زیاد می‌باشد که این امر نیز مدل کردن حرکتها را دشوار ساخته است. در نتیجه تولید واسطه‌هایی که بتوانند حرکتهای زیادی را با هر درجه پیچیدگی بصورت بی‌درنگ^۳ فهمیده و دستورات معادل آنها را اجرا کنند، کاری دشوار است. چند نمونه از حرکتهایی که برای آنها مدلها و الگوریتمهایی ارائه شده است عبارتند از: حرکت دست، مدل کردن حرکت سر و صورت، حرکت اندامهای صورت مانند لبها، ابروها و چشمها که در ادامه استفاده از این حرکتها را بطور مختصر بررسی می‌کنیم.

استفاده از حرکات دست بعد از صحبت کردن بیشترین کاربرد را در برقراری ارتباط بین انسانها دارد و از گذشته‌های دور به کسانی که توانایی سخن گفتن را نداشتند حرکات ایما و اشاره‌ای آموزش داده می‌شد که بتوانند با حرکتهای مختلفی که با دست انجام می‌دهند، منظور خود را برسانند. حتی انسانهایی که توانایی سخن گفتن را نیز دارند در هنگام صحبت کردن و تعامل با سایر انسانها از حرکات دست استفاده می‌کنند. بطور مثال ما برای خواستن چند شیء، چند انگشت دست خود را نشان می‌دهیم و یا برای اینکه توجه شخصی را به چیزی جلب کنیم، با انگشت دست خود به آن اشاره می‌کنیم. در حقیقت اشاره کردن با دست آنقدر در بین انسانها معمول است که بعضی از آنها حتی وقتی با تلفن صحبت می‌کنند نیز با دست خود حرکات اشاره‌ای انجام می‌دهند. در یک واسطه بینایی مبتنی بر حرکت دست، معمولاً سه دسته اطلاعات از تصاویر استخراج می‌شود. اول اینکه در هر تصویر مکان دست در تصویر تشخیص داده می‌شود و سپس از روی مکان دست در فریمهای متوالی تغییرات مکانی دست مشخص می‌شود. همچنین در

^۱ Interface

^۲ VBI (vision-based interaction)

^۳ Real-time

بسیاری از این سامانه ها (سیستم ها) هر وضعیت^۱ دست برای خود یک معنا و مفهوم جداگانه ای دارد و بنابراین تشخیص نوع حرکت نیز یکی از عملیاتهای مورد نیاز برای این واسط ها می باشد. در بعضی از الگوریتمها یک سری اطلاعات پارامتریک دیگر مانند جهتی که انگشت دست اشاره می کند نیز استخراج می شود.

در استفاده از حرکت لبها در ابتدا باید به این موضوع اشاره کنیم که تشخیص گفتار^۲ یکی دیگر از روشهای تعامل انسان با کامپیوتر (HCI)^۳ می باشد که با استفاده از این روش، کامپیوتر قادر خواهد بود که صحبتهای کاربران را شنیده و اطلاعات گفتاری آنها را استخراج کند. در حال حاضر بسیاری از سامانه های تشخیص گفتار موجود، تنها در محیط آزمایشگاهی عملکرد بسیار خوبی دارند و در محیطهای واقعی بدلیل نویز بالای محیط، این سامانه ها بخوبی عمل نکرده و کارایی آنها کاسته می شود. در مقابل این سامانه ها، انسانها قادرند که با ترکیب سیگنالهای گفتار و اطلاعات بصری دیگر مانند حرکت بدن و صورت، در محیطهای نویزی نیز تشخیص گفتار را بصورت قابل قبولی انجام دهند که این امر در بین افرادی که توانایی شنیدن را ندارند بوضوح دیده می شود. با الهام گرفتن از نحوه عملکرد انسان در تشخیص گفتار و اینکه اطلاعات بصری گفتار در محیطهای نویزی تغییر نکرده و حساس به نویزهای صوتی نیستند، این انگیزه در بین محققان بوجود آمد که از اطلاعات بصری نیز در سامانه های تشخیص گفتار استفاده کنند. از آنجا که سیگنالهای صوتی از مجرای گفتار انسان خارج می شوند، در تمامی سامانه های تشخیص گفتار مبتنی بر اطلاعات بصری (سامانه های صوتی- بصری)^۴ سعی شده است که حرکت دهان و لبها مدل شود و اطلاعات استخراج شده از این حرکتها با اطلاعات صوتی ترکیب شده و یک سامانه مقاوم در مقابل نویز ایجاد شود. در واقع در اینگونه سامانه ها، تعامل بصورت ترکیبی از بینایی ماشین و سیگنالهای صوتی می باشد.

یکی دیگر از حرکتهایی که در واسطهای انسان - کامپیوتر مورد استفاده قرار گرفته است، حرکت اندامهای صورت بخصوص چشمها و ابروها می باشد. اگر چه استفاده از حرکات دست بیشترین کاربرد را در واسطهای مبتنی بر بینایی دارد، اما اخیرا محققان در پی این بوده اند که با ابداع روشهای جدید، واسطهای موجود را طبیعی تر کرده و هر چه بیشتر به شیوه تعامل انسان- انسان نزدیک تر کنند. یکی از این روشها استفاده از حرکت چشم برای انتخاب یک منو بر روی صفحه نمایش می باشد. انسانها برای اینکه کاری را

^۱ Posture

^۲ Speech Recognition

^۳ Human-Computer Interaction (HCI)

^۴ Audio-Visual

روی یک شی انجام دهند ابتدا به شی نگاه کرده و سپس با استفاده از حرکت دست، کاری را روی شی انجام می‌دهند. همچنین هنگامی که یک کاربر از کامپیوتر استفاده می‌کند، معمولاً ابتدا به قسمت مشخصی از صفحه نمایش نگاه کرده و سپس کار خود را انجام می‌دهد. در بعضی از واسطه‌های مبتنی بر حرکت دست، کاربر ابتدا با نگاه کردن شی مورد نظر خود را پیدا کرده و سپس با حرکت دست کاری را روی آن انجام می‌دهد و نگاه کردن به شی قبل از و سریعتر از حرکت دست انجام می‌شود، مانند محیط‌هایی که برای واقعیت‌های مجازی^۱ طراحی شده‌اند. همچنین برای معلولینی که توانایی حرکت دادن دستهای خود را ندارند نگاه کردن راحت‌ترین کار ممکن برای استفاده از کامپیوتر می‌باشد. روش کار با واسطه‌های مبتنی بر حرکت چشم بدین صورت است که روی صفحه نمایش منوهای سامانه گنجانده می‌شود و کاربر در زمانی حدود چند ثانیه به این منوها خیره شده و سپس سامانه منوی مورد نظر کاربر را انتخاب می‌کند. همچنین چند الگوریتم ترکیبی نیز ارائه شده است که در این روشها با استفاده از نگاه کردن و یک ورودی دیگر مانند حرکت دست، حرکت ابروها و یا استفاده از موشواره، کاربر سریعتر کارهای خود را انجام می‌دهد. در چند نمونه از این سامانه‌ها محققان با استفاده از روش برق‌نگاری ماهیچه (EMG)^۲ تغییر سیگنالهای الکتریکی حرکات صورت را اندازه گرفته و با تلفیق این روش با حرکت چشم، واسطه‌هایی را طراحی کردند. ایده این موضوع از روی تعامل‌های انسان با انسان گرفته شده است، بدین صورت که در بسیاری از موارد انسانها با استفاده از حرکتهای قسمت‌های مختلف صورت مطلبی را بیان کرده و یا حالتی را نشان می‌دهند. در یکی از این واسطه‌ها با استفاده از روش EMG تغییر سیگنالهای الکتریکی ناشی از حرکت ابرو با داده‌های خیره شدن چشم ترکیب شده و کاربر بعد از اینکه به یک شی خاص خیره می‌شود با حرکت دادن ابروهای خود آن شی را انتخاب می‌کند [۱].

از حرکتهای سر و صورت نیز در سامانه‌های HCI استفاده شده است. در این سامانه‌ها کاربر می‌تواند با حرکت دادن سر و صورت در جهت‌های مختلف با سامانه ارتباط برقرار کند و دستوراتی را به سامانه منتقل نماید. یکی از مشکلات این سامانه‌ها این است که تعداد حرکتهای مختلفی که می‌توان با سر و صورت انجام داد محدود است و کاربران نیز تمایل زیادی به استفاده از حرکت سر و صورت برای تعامل ندارند، ولی می‌توان این حرکتهای را بصورت ترکیبی با سایر حرکتهای بکار برد.

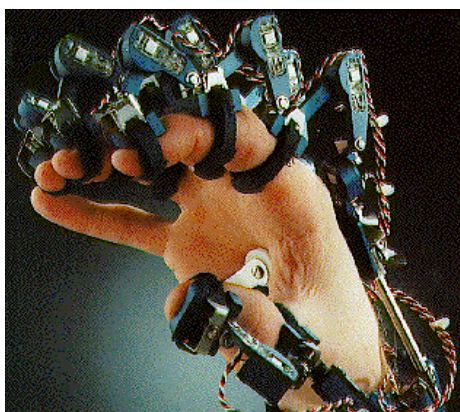
از بین تمامی حرکتهایی که در این بخش به آن اشاره شد، حرکت دست بیشترین کاربرد را در واسطه‌های انسان - کامپیوتر داشته است و بیشترین الگوریتمها نیز در مورد حرکت دست ارائه شده است.

^۱ Virtual Reality

^۲ Electromyography (EMG)

علت اینکه حرکت دست بیشترین کاربرد را دارد یکی این است که تعداد حرکت‌هایی را که می‌توان با دست انجام داد از اندام‌های دیگر بدن بیشتر است و این تنوع حرکتها استفاده از سامانه‌ها را آسانتر می‌کند و اجازه می‌دهد که قابلیت‌های بیشتری در سامانه گنجانده شود. دلیل دوم این است که یادگیری حرکت‌های دست برای کاربران راحت‌تر بوده و حتی کاربران غیر حرفه‌ای نیز می‌توانند در زمان کوتاهی تمامی حرکتها را یاد بگیرند که این موضوع باعث شده است که کاربران تمایل بیشتری برای استفاده از حرکت‌های دست نشان دهند. در ادامه استفاده از حرکت‌های دست در سامانه‌های HCI را بررسی می‌کنیم.

همانطور که پیشتر نیز اشاره شد، در چند دهه اخیر استفاده از حرکت‌های دست در سامانه‌های HCI بسیار مورد توجه قرار گرفته است. در ابتدا برای اینکه کامپیوتر (ماشین) بتواند حرکت‌های دست را تفسیر کند از یک سری دستگاه‌های الکترومغناطیسی استفاده می‌شد که می‌توانستند مکان دست و بازوها و زوایای بین آنها و حتی زوایای بین مفاصلهای انگشتان را اندازه بگیرند. به این گروه از دستگاهها، دستکش‌های داده‌ای^۱ و به این روش برقراری ارتباط بین انسان و کامپیوتر، روش مبتنی بر دستکش^۲ گفته می‌شود. دلیل این نامگذاری این است که اینگونه دستگاهها بصورت دستکش‌هایی هستند که سنسورهایی روی آن قرار دارد و با تعداد زیادی سیم به کامپیوتر متصل می‌شود و کاربران برای استفاده از سامانه باید آن را مانند یک دستکش بپوشند. استفاده از اینگونه دستکشها چندان مورد استقبال کاربران قرار نگرفت، چرا که استفاده از ابزارهایی نظیر موشواره و صفحه کلید بسیار ساده‌تر از این دستکشها بود و در حال حاضر از این دستکش‌ها تنها برای کارهای خاص و کنترل در بعضی محیط‌ها استفاده می‌شود. مثلاً چنین روشی را برای شبیه‌سازی جراحی در یک محیط مجازی می‌توان به کار برد. در شکل ۱-۱ یک نمونه از این دستکشها نشان داده شده است.

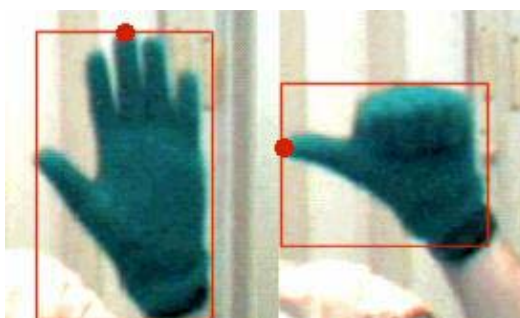


شکل ۱-۱ یک نمونه دستکش داده‌ای

^۱ Data Gloves

^۲ Glove-Based

برای برطرف کردن محدودیتهای دستکشهای داده‌ای، استفاده از بینایی ماشین برای تشخیص حرکت‌های دست پیشنهاد شد. در این رویکرد پیشنهاد شد که با استفاده از دوربینهای ویدئویی و روشهای بینایی ماشین و پردازش تصویر، حرکت‌های دست شناسایی و تفسیر شود. در ابتدا به جای استفاده از دستکشهای داده‌ای، پیشنهاد شد که برای متمایز کردن دستها از پس زمینه^۱ از دستکشهای رنگی استفاده شود و یا علامت‌های رنگی بر روی دست کاربران قرار داده شود. یک نمونه از این دستکشها در شکل ۱-۲ نشان داده شده است.



شکل ۱-۲ استفاده از دستکشهای رنگی به جای دستکشهای داده‌ای

استفاده از این دستکشها در مقایسه با دستکشهای داده‌ای بسیار مناسب‌تر می‌باشد ولی همین دستکشها نیز یک محدودیت برای سامانه‌های HCI به شمار می‌رود، چرا که کاربر برای استفاده از سامانه همیشه مجبور است از این دستکشها استفاده کند. برای حل این مشکل طراحان از سامانه بینایی انسان الهام گرفتند. ما انسانها می‌توانیم تنها با استفاده از شکل و رنگ دستها، مکان دست در تصویر را مشخص کنیم. این موضوع این ایده را در ذهن طراحان بوجود آورد که به جای استفاده از دستکشهای داده‌ای و رنگی، تنها از شکل دستها و رنگ پوست برای تشخیص مکان دست در تصویرها استفاده کنند.

در تحقیقات انجام شده از دو نوع حرکت دست در تعامل انسان با کامپیوتر استفاده می‌شود. در بعضی از سامانه‌ها ترکیب چند وضعیت^۲ دست در فریمهای متوالی معادل یک دستور در نظر گرفته می‌شود، حال آنکه در دیگر سامانه‌ها، وضعیت دست در هر فریم معنا و مفهوم مستقلی دارد. به حرکت‌های نوع اول حرکت‌های پویا^۳ و به نوع دوم حرکت‌های ایستا^۴ گفته می‌شود. در حرکت‌های ایستا تنها باید شکل دستها

^۱ Background

^۲ Posture

^۳ Dynamic Gesture

^۴ Static Gesture

در هر فریم شناسایی شود، حال آنکه برای حرکت‌های پویا باید ردیابی دست‌ها نیز انجام شود تا علاوه بر شکل دست‌ها در فریم‌های متوالی، نحوه جابجایی آنها نیز مشخص شود.

تاکنون تحقیقات زیادی برای استفاده از حرکت‌های دست انجام شده است. در [۲] آقای گوپتا الگوریتمی برای تشخیص ۱۰ حرف از زبان اشاره‌ای آمریکایی^۱ ارائه داده است. این زبان یکی از قویترین زبان‌های اشاره‌ای می‌باشد که افراد ناشنوا می‌توانند برای بیان کلمات و جملات انگلیسی از آن استفاده کنند. در این الگوریتم از دست کاربران در زمینه سیاه تصویر گرفته می‌شود. پس از تصویر برداری، با استفاده از روش آستانه‌سازی هیستوگرام، تصویر دست از بقیه قسمت‌های زمینه تقطیع می‌شود. آنگاه با استفاده از روش توالی پیرامونی محلی، بردارهای ویژگی از پیرامون دست‌ها استخراج شده و با روش تطبیق خطی، وضعیت‌های دست شناسایی می‌شود. در این الگوریتم تصویر برداری از حرکات دست در محیط آزمایشگاهی با زمینه ثابت و نور فلورسنت از بالا انجام شده است و شناسایی هر حرکت حدود ۲۰ ثانیه طول می‌کشد. در این الگوریتم تنها از حرکت‌های ایستا استفاده شده است.

در [۳] آقای الجراح برای شناسایی ۳۰ حرف زبان اشاره‌ای عربی از شبکه‌های عصبی-فازی استفاده کرده است. در این الگوریتم نیز از تصاویر سطح خاکستری در زمینه ثابت استفاده شده است. برای تقطیع دست‌ها از یک الگوریتم آستانه‌سازی تکراری استفاده می‌شود. برای استخراج ویژگی ابتدا مرکز ثقل دست‌ها محاسبه شده و سپس بردارهایی از این نقطه به نقاط پیرامون وصل می‌شود. طول این بردارها مولفه‌های بردار ویژگی را تشکیل می‌دهند. بردارهای ویژگی به ۳۰ شبکه داده شده و شبکه‌ای که خروجی بزرگتری تولید می‌کند، نوع حرکت دست را مشخص می‌نماید. در این الگوریتم نیز تنها از حرکت‌های ایستا استفاده شده است.

در [۴] آقای تنگ الگوریتمی بلادرنگ برای شناسایی ۳۰ حرف زبان اشاره‌ای چینی ارائه کرده است. در این سیستم از هیچ نشانه‌ای بر روی دست‌ها استفاده نمی‌شود و زمینه تصاویر نیز ساده نیست و همچنین هیچ نورپردازی خاصی برای سیستم در نظر گرفته نشده است. برای تقطیع تصاویر دست، از روش تفاضل فریم‌های متوالی و یک مدل رنگ پوست استفاده شده است. در این الگوریتم کل تصویر دست به عنوان بردار ویژگی در نظر گرفته می‌شود و برای شناسایی حرکت‌ها از روش نزدیکترین همسایگی استفاده شده است. در این الگوریتم نیز تنها از حرکت‌های ایستا استفاده شده است.

^۱ American Sign Language (ASL)

در [۵] آقای لیکسار و همکارش یک الگوریتم تطبیقی برای شناسایی ۹ وضعیت ایستا ارائه داده است. تاکید نویسندگان این مقاله بر ارائه سیستمی بوده است که بتواند حتی با تعداد کم نمونه های آموزشی نیز قدرت تعمیم خوبی داشته باشد و با دقت مناسبی حرکات کاربران جدید را تشخیص دهد. الگوریتم ارائه شده در این مقاله این قابلیت را دارد که بصورت برخط توسط کاربران آموزش مجدد داده شود. این آموزش می تواند در دو حالت با نظارت و بی نظارت انجام شود. حالت با نظارت زمانی فعال می شود که حرکتی اشتباه تشخیص داده شود، ولی زمانیکه سیستم حرکتهای را درست تشخیص می دهد نیز آموزش انجام می شود (یادگیری بی نظارت). در حالت بی نظارت، اگر سیستم حرکتهای کاربر را درست تشخیص دهد، پارامترهای این حرکتهای را جایگزین نمونه های ذخیره شده می کند. آزمایشات انجام شده در مقاله بدین صورت است که یک تصویر با استفاده از ویدئو پروژکتور بر روی دیوار نشان داده می شود و کاربر بصورت بلادرنگ می تواند اسلاید قبلی یا بعدی را نمایش دهد و یا اینکه بر روی اسلاید جاری تغییراتی اعمال کند. برای تقطیع دستها از روش تفاضل پس زمینه استفاده شده و تصویر پس زمینه همان تصویری است که ویدئو پروژکتور بر روی دیوار منعکس می کند. برای بردارهای ویژگی از ۶ مولفه اول توصیفگرهای فوریه استفاده شده و دسته بندی حرکتهای با استفاده از روش نزدیکترین همسایگی انجام می شود.

آقای چانوا و همکارش در سال ۲۰۰۲ یک واسط گرافیکی ساده پیاده سازی کردند که کاربر می تواند پنجره ها و اشیاء موجود در آن را تنها با استفاده از حرکات دست انتخاب کرده، جابجا کند و یا اندازه آنها را تغییر دهد [۶]. همچنین کاربر می تواند پنجره ها را باز و بسته نماید و یا کارهای انجام شده را به حالت اول برگرداند. برای این واسط گرافیکی ۱۴ حرکت در نظر گرفته شده که از بین این حرکتهای چهار حرکت ایستا و مابقی پویا هستند. همچنین شش حرکت این واسط تنها با یک دست انجام داده می شوند و برای مابقی، کاربر باید از هر دو دست خود استفاده کند. برای تقطیع تصویر دستها از پس زمینه، دو روش تفاضل پس زمینه و مدل رنگ پوست استفاده شده است. در این واسط محدودیتی برای پوشش کاربر در نظر گرفته نشده و از آنجا که کاربر می تواند از پیراهن های آستین کوتاه نیز استفاده کند، طراحان سیستم برای حذف قسمتهای ناخواسته دست از روشهای برش مچ استفاده کردند. برای استخراج ویژگیها از پیرامون دستها و تبدیل فوریه گسسته استفاده شده است. شناسایی دستها با کمک شبکه های عصبی RBF و مدل مخفی مارکوف انجام می شود و این واسط می تواند با دقت ۹۲ درصد و بصورت بی درنگ حرکتهای را تشخیص داده و معادل دستوری آنها را اجرا کند.

در [۷] آقای یوانسین ژو و همکارش یک مرورگر نقشه^۱ پیاده سازی کردند که کاربر می تواند با استفاده از ۱۲ حرکت پویا تصاویر را به بالا، پایین، چپ و راست جابجا کند، بزرگ نمایی یا کوچک نمایی کرده و حتی در زوایای مختلف دوران دهد. کلیه حرکت های این واسط با یک دست انجام می شود. برای تقطیع تصاویر دست از روش تفاضل فریم های متوالی و یک جدول جستجو^۲ برای تشخیص نواحی به رنگ پوست استفاده شده است. از آنجا که حرکت های این مرورگر نقشه حرکت های پویا هستند، لازم است برای تشخیص حرکات هم اطلاعات مربوط به شکل دستها و هم اطلاعات مربوط به جابجایی دستها استخراج شود. برای استخراج ویژگی های حرکتی از روش رگرسیون حرکت تصویر^۳ استفاده شده است. با استفاده از این روش، ویژگی های مربوط به حرکت های دست مشخص می شود. برای استخراج ویژگی های مربوط به وضعیت دستها از گشتاورهای هندسی مرتبه اول و دوم استفاده شده است. ترکیب ویژگی های حرکتی و ویژگی های مربوط به وضعیت دستها در هر فریم، بردارهای ویژگی را تشکیل می دهند. برای شناسایی حرکتها از روشی مشابه روش نزدیکترین همسایگی استفاده شده است. این واسط می تواند در حالت برخط با دقتی بین ۸۰ تا ۹۰ درصد و در حالت برون خط با دقت ۹۰/۸۳ درصد حرکتها را شناسایی کند. کل زمان لازم برای تقطیع تصاویر، استخراج ویژگی و شناسایی حرکتها حدود ۲ ثانیه می باشد.

در [۸] آقای استارنر و همکارانش دو سامانه برای شناسایی ۴۰ علامت زبان اشاره ای آمریکایی پیاده سازی کرده اند. در سامانه اول دوربین بر روی یک میز قرار دارد و کاربر در مقابل دوربین حرکت های دست را انجام می دهد. در سامانه دوم یک دوربین بر روی کلاه کاربر نصب شده است که از بالا حرکت های دست کاربر را تصویربرداری می کند. برای تقطیع تصاویر دست از مدل رنگ پوست و روش رشد دادن نواحی استفاده شده است. پس از تقطیع تصاویر دست ها، از گشتاورهای درجه دوم برای استخراج ویژگی استفاده شده و ۱۶ ویژگی از هر دست استخراج می شود. در حالتی که همپوشانی دستها رخ می دهد تنها یک ناحیه به رنگ پوست استخراج می شود. در این حالت ویژگی های استخراج شده این ناحیه به هر دو دست تعلق می گیرد. علاوه بر این در نمونه های آموزشی نیز چند مورد همپوشانی در نظر گرفته شده است تا دقت سامانه ها بخاطر همپوشانی دستها کاهش نیابد. برای شناسایی حرکتها نیز از مدل مخفی مارکوف استفاده شده است. برای سامانه اول که دوربین بر روی میز قرار داده شده است، دقت ۹۲ درصد و برای سامانه دوم که دوربین بر روی کلاه کاربر تعبیه شده است، دقت ۹۷/۸ درصد می باشد. دلیل این برتری در این است که در سامانه دوم مشکل همپوشانی دستها کمتر وجود دارد.

^۱ Map Browser

^۲ Look up Table

^۳ Image Motion Regretion