Paper ID 1106

#### **Real-time Gesture Recognition for Interactive Presentation with Depth Sensor** Anonymous CVPR submission Abstract In this paper, we propose a method to enable real-time interaction between the projection contents and speaker through detecting and recognizing meaningful human ges-tures from depth maps captured by depth sensor, making projection screen as a kind of touch screen. Considering that depth noise and serious occlusion may ruin the con-struction of skeleton, our hand trajectory is derived from Potential Active Region. To cope with their inter-class and intra-class variations, hand trajectory is temporally seg-mented into movements, which are represented as Motion History Images. A novel set-based soft discriminative mod-el is learned to recognize gestures from these movements. In addition, as it is a real-time system, a complexity reduction method is employed. The proposed approach is evaluated on our dataset and performs efficiently and robustly with 90% correct recognition rate. **1. Introduction** Recently, gesture recognition has been attracting a great deal of attention as a natural human computer interface, s-ince it allows users to control or manipulate devices in a more natural manner through intentional physical move-ments of figures, hands, arms, face, head, or body. So far, numerous studies [10][16] have been conducted on gesture

recognition for human computer interaction, especially for hand and arm gesture recognition. Based on these technologies, a number of recognition applications are developed, including sign language recognition, game controlling, navigating in virtual environment, etc. In this paper, we present an effective and efficient arm gesture recognition algorithm that enable natural interaction between speaker and presentation contents. During presentation, speakers often s-tand in front of the projection screen at a distance from the machine with projection contents. It is more natural for s-peakers to remote control the page flipping, scrolling and clicking by arm gestures. Considering the potential light influences caused by projector on the color images, depth maps captured during presentation are employed for ges-



Figure 1. Framework of proposed approach.

ture recognition in our work. Fig. 1 shows the framework of proposed approach. In the framework, human body is segmented from noisy depth maps, and then potential active regions (PAR) are derived from head position for meaningful gesture detection. Once the hand is observed in the potential active region, its trajectory will be recorded and decomposed to a series of movements. (see green box in dash line). These movements are represented as Motion History Images (MHI) and assembled to a labeled gesture by utilizing proposed set-based soft discriminative model.

As depth data is noisy, background suppression technique is used to segment human body from background. The location and size of human body are determined by searching a proper bounding box in the generated human body's depth map. Considering that human body may be incomplete in that bounding box because of occlusion or limited view angle of depth sensor, head detection is applied for estimating the size of the complete human body. However, normal method using face to detect head is impossible in presentation since speakers may turn their faces to the

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

Up				And a second sec		A BAC
Down			And A State			And C
Go						
Back			Manual Andrew Construction of the second sec		A the first of the	
Click	1	ř j	-	-		2

Figure 2. Five gestures ("click" is shown in depth since it is a gesture vertical to the plane).

side (not facing the sensor) and skin color is also changed due to the strong light from the projector. As a result, the only useful cue to detect head is the depth map.

Potential Active Regions of human arms (see the boxes beside body in Fig. 1) are adaptively determined by the location and size of body. Intuitively, PARs are the most discriminative regions. Therefore arm is only detected when reaching in PARs, which is shown in the framework of our system in Fig. 1. If no arm is detected in PARs, it is unnecessary to perform recognition step in this frame or even unnecessary to detect arm in the next few frames. That vastly reduces the computational complexity. Once arm is observed in a PAR, the trajectory of hand will be recorded and decomposed to a series of movements (see the green box in dash line), which are classified by Support Vector Ma-138 chine (SVM). Offline training is not required to decompose 139 the trajectory. These movements are assembled into a set, 140 which is then labeled and understood as a specific gesture. 141 Considering that one misclassified movement in the set may 142 influence the result of gesture recognition. Original propos-143 144 al of a set-based soft discriminative model can correct the misclassified movement and designate a most likely gesture 145 146 label to the set. Experimental results show that this model has a better performance than traditional one. 147

Our main contribution lies in three folds. First, the noise 148 and occlusion problems caused by depth sensor are solved 149 by efficient preprocessing. Second, with PAR, the compu-150 151 tational complexity is dramatically reduced by selecting active frames. Third, a set-based soft discriminate model is 152 153 originally proposed, which has the ability to correct misclassified movements. To test our method, we capture a 154 dataset including 5 gestures: "up", "down", "go", "back" 155 and "click" (see Fig. 2). 156

The rest of this paper is organized as follows. Related work is introduced in section 2. Section 3 presents the
proposed method and the experimental results in section 4
demonstrate the efficiency of the proposed method. Section
5 concludes this paper and gives the future directions.

## 2. Related Work

CVPR 2013 Submission #1106. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

In the last decade, most of recognition algorithms are designed for the color images captured by monocular camera sensors. One challenge of the color image based recognition is how to efficiently segment the object from the background. In order to obtain foreground silhouettes of objects, most of the recognition algorithms are restricted to pure and static backgrounds. For example, public dataset KTH human motion dataset [12] and Weizmann human action dataset [1] both record human actions under relatively static backgrounds. The rapid development of depth sensors open up the possibilities of dealing with cluttered background by providing depth information. Even though, depth sensors like Kinect [14], Time of Flight camera [5] or stereo camera [15] still present two challenges: noise and occlusion. Noise decreases the quality of background subtraction mainly because of the limited measurement accuracy of the depth sensor. Occlusion occurs when there is an object (e.g. desk) in front of human body. It will ruin body detection because part of the human body is blocked. This is also a main problem when adapting monocular camera. What is worse, Kinect sensor regards black objects, such as black trousers or black hair, as occlusions due to its generation principle. To overcome the occlusion problem, some approaches employ the location of human face to indicate the location of human body. Face detection is usually implemented as the first step to detect human body. For example, Wang et al. [18] locate the face before obtaining skin model from face and use it to detect human hand. As stated before, they cannot work well in presentation due to the variant directions of face and abnormal skin color, which is influenced by projector light.

When equipped with depth sensor, many researchers make effort to compute 3D joint positions of human skeleton. Shotton et al. [14] provide a rather powerful human motion capturing technique. There are also many action recognition works start directly from human skeleton. For example, Jiang et al. [17] track 20 joint positions by the skeleton tracker proposed by [14] and use local occupancy pattern to represent the interaction. Sadeghipour et al. [11] also directly use the 3D joint positions of human skeleton for gesture-based object recognition. Both of them use Kinect as the depth sensor. When referred to the other depth sensors, robust and fast method has not been presented yet. That means, methods in [17] and [11] are restricted to Kinect sensor. Besides, distance from human and Kinect in both of their datasets are fixed while in a normal presentation, speakers walk around. Therefore, to make our system more natural for interaction and more general for other stereo sensors, skeleton are replaced with silhouette. By accumulating the silhouettes in PARs, MHIs are generated.

PAR is a spatial region (where), the generation of MHI still asks for a temporal region (when). In another word, the

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

216 start and end of a gesture in a video sequence should also 217 be determined. Videos in most published 3D datasets are 218 readily segmented into sequences that contain one instance 219 of a known set of action labels. For example, Sadeghipour 220 et al. [11] capture the 3D Iconic Gesture dataset and seg-221 ment the video by the moment when subject retracting their 222 hands back to the rest position, so does NATOPS Aircraft 223 Handling Signals Database captured by Song et al. [15]. 224 Different from their works, trajectory decomposing is em-225 ployed to automatically segment the video in our approach 226 to provide a more natural interaction.

227 As we know, representation of suitable feature and mod-228 eling of dynamic patterns are the two main important is-229 sues. In our work, MHI serves as the representation of the 230 movement in the action recognition framework according to 231 the taxonomy summarized by Bobick [2] in an early survey. 232 Similarly, 3D low-level features are deeply studied in recent 233 years. Most of them are extended from normal 2D features 234 such as [9] (3D Harris corner detector), [20] (3D SURF de-235 scriptor), [8] (3D HOG descriptor) and [13](3D SIFT de-236 scriptor). However, these local featurLes are not discrimi-237 native features in textureless depth maps. To assemble low-238 level movements into a gesture in the proposed framework, 239 generative model and discriminative model are the main t-240 wo temporal state-space models. Generative model learns 241 to model each class individually and always assume that the 242 observations in time are independent, e.g., Feng et al. [4] 243 and Weinland et al. [19]. Discriminative model is trained 244 to discriminate between action classes and model a condi-245 tional distribution over action labels given the observation. 246 Jordan et al. [7] compare discriminative and generative 247 learning as typified by logistic regression and naive Bayes. 248 Our gesture recognition model belongs to the discriminative 249 model. 250

## 3. Method

251

252

253

254

255

256

257

258

259

260

261

262

263

In our framework, the background is firstly suppressed from noisy depth maps, and then PARs are derived from head position for meaningful gesture detection. The generation of PARs and the segmentation of the hand trajectory is to determine where and when to generate MHI in the video sequence. Then, a soft discriminative model is originally proposed to assemble them into one gesture. In addition, with hand detecting in the PARs, the complexity of the system is vastly reduced.

## 3.1. Background Suppression

In presentation, touching the screen while performing gestures is a natural way to control the projection contents. However, measurement accuracy of the depth sensor is limited. It is difficult to segment objects from the background when the depth of objects is too close to that of background. As the result, even if the depth of background has been cap-



Figure 3. (a) Results of improved background suppression method. (b) Method to obtain the size and location of head. (c) Potential active regions for two arms. (d) Distance map using chamfer distance.

tured, segmenting arms from screen is impossible by background subtraction techniques. It is worth noticing that the depth value of each pixel in the background is approximately Gaussian distributed with time going on, and the depth value at a fixed position lies in a certain range with a high probability. When the probability of the depth value is below a threshold, the corresponding point is processed as human body. Fig. 3 (a) shows the improvement.

## 3.2. PAR Generation and Hand Trajectory Decomposition

The size and location of the whole human body is necessary for PAR generation. Therefore, a bounding box containing the whole body is constructed. The segmented depth map is scanned along vertical lines from left to right while recording the proportion of body pixels in the line. Location of the left border of bounding box is determined when the first time the proportion reaches a threshold. Locations of other three borders are searched in a similar way.

Notice that human body in such bounding box may be 307 incomplete because of the occlusion or the limited view an-308 gle of sensor. Size of head is thus employed to estimate the 309 size of the whole body according to the normal proportion 310 of human figure. As we know, physical width of human 311 body varies with height. In generally, the width of neck is 312 the smallest while that of shoulder is the largest. Inspired by 313 this observation, pixels' values are summed along horizon-314 tal direction in the body bounding box. A curve is drawn 315 to show the proportions of human pixel in each horizontal 316 line (see Fig. 3 (b)). After smoothing the curve by media 317 filter, location and size of head are obtained by detecting the 318 largest width variation. Two PARs are then constructed for 319 left and right arms. Since PARs cover the most likely arm 320 movement region, they can serve as the constraint for hand 321 position tracking and adaptive detection mode transfer. The 322 detection mode transition will be introduced in details in the 323

CVPR #1106

### CVPR 2013 Submission #1106. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431



Figure 4. Decomposing the movement trajectories. The start and end points in each segmentation are interest points. 10 movements are detected in this example.

section 3.3.

A chamfer distance map is derived when hand enters the PAR. As the farthest endpoint from the edge of PAR, hand is detected and tracked in that region before leaving (see Fig. 3 (d)). However, it is difficult to recognize gestures directly from the trajectories of hand since our dataset has the following two challenges. First, the dataset has small inter-class variances. Gesture "up" is very similar to "down" while the trajectories of "go" and "back" are similar if not considering the order of frame. Second, it also has large intra-class variances. There are variant ways to act the same gesture by different subjects. Even the same subject performs the same gesture differently each time. To recog-350 nize gestures from the complex trajectory, a method is pro-351 posed to decompose the trajectory into movements, classify 352 the movements and then assemble them to form a mean-353 ingful gesture. See from Fig. 4, the subject successively 354 performs four actions in one video. Among the complex 355 trajectory, interest points are required to help segmenting 356 it into movements. As we know, interest point is the sud-357 den change of trajectories in video sequence. Based on that 358 common sense, the method to search the interest points is 359 introduced as follows. 360

Method of Least Squares (MLS) is used to detect the in-361 terest points on the trajectories. The MLS assumes that the 362 best-fit curve of a given type is the curve that has the mini-363 mal sum of the deviations squared (least square error) from 364 a given set of data. The type of straight line y = ax + b is 365 employed to approximate a given set of points. A new point 366 is determined as an interest point when deviation  $d_{n+1}$  of 367 a new point  $(x_{n+1}, y_{n+1})$  is larger than a threshold  $d_{thre}$ , 368 otherwise not. Our method avoids complex off-line train-369 ing. So long as an interest point is found, the trajectory 370 is segmented, which indicates that a movement is detected 371 (see Fig. 4). 372

## **3.3. Detection Mode Transfer**

373

374

Except defined gestures, most of the arm movements of
the speaker are meaningless to the system. If hand trajectories are far away from screen or not in the PARs, it is unnec-



Figure 5. (1) Hand is detected in PARs. (2) Hand leaves PARs. (3) Hand leaves PARs for a long time.

essary to record the hand trajectory or recognize the gesture. The potential hand position can be predicted from the hand position in the previous frame. The detecting frequency can be reduced if the hand position is far from the PAR or its depth out of defined depth range. To adjust the detection frequency to arm movements, three hand detection modes are defined as "inactive", "semi-active" and "active", and each mode has different detection intervals k. In an inactive mode, hands will be re-detected after K frames, that is, k = K - 1. In the active mode, hands will be re-detected in the next frame, i.e., k = 0. If the detection mode keeps semi-active, the interval k is linearly increases from 1 to (K - 1).

When a hand is detected in PARs and is close enough to screen, the mode is immediately switched to "active" from "inactive" or "semi-active" mode. Otherwise, the mode is switched to "semi-active" mode. Since the number of frames being detected decreases in "semi-active mode". It will finally become "inactive mode" (see Fig. 5). It should be noted that "active" mode is not directly changed to "inactive" mode. Because the detecting module may generate a false reject error, i.e., missing a hand in one frame. In "Semi-active" mode, hand can be re-detected after less frames than in "inactive" mode. Notice that changing mode from "inactive" to "active" may require a relatively long time, during which hand may have already entered the PARs in those skipped frames. To avoid this phenomenon, an extension of PAR is generated for pre-detecting hand to ensure the completeness of trajectory. The hand trajectory does not include the hand detected in extension of PAR.

### 3.4. Feature extraction and classification

Before feature extraction, the size of PAR is normalized to facilitate the generation of the uniform MHI, which is used to represent decomposed movements. We adopt two stages classification pipeline to extract global feature, i.e., coding and pooling. The model is trained and tested with this feature by SVM.

**Coding.** The silhouettes of human arm belonging to the same movement in PAR are accumulated to generate the MHI, which is originally proposed by Bobick et al. [3]. In a MHI, pixel intensity records the temporal history of motion

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451 452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539



Figure 6. (a) Samples of 12 classes of movements. Class 12 includes meaningless movements. Three of them are listed. (b) Uniform rectangle spatial region. (c) Uniform semi-circle spatial region.

at each position. The MHI of decomposed movements are shown in Fig.6 (a).

**Pooling.** Designing a proper spatial region for pooling has a significant impact on feature's distinguishing capacity [6]. Since arms rotate along shoulder joints, the pixel intensity distributes in a semi-circle manner. Feature pooled from the uniform spatial region like that in Fig. 6 (b) is improper for further classification. For more effective representation, we choose the spatial regions shown in Fig. 6 (c). The blocks in the semi-circle region are represented as  $R_1, R_2, ..., R_N$ . The final global feature is denoted as  $V = v_1, v_2, ..., v_N$  and  $I_{m,n}$  is the intensity of pixels at the position (m, n) in the MHI.  $v_i$  is computed as follows,

$$v_i = \frac{1}{|R_i|} \sum_{(m,n)\in R_i} I_{m,n}, i = 1, 2, \dots, N$$
 (1)

where  $|R_i|$  represents the pixel number in this block.

**Classification.** SVM is applied for classifying movements. The standard SVM divides two classes by clear gap that is as wide as possible. The classifier predicts new examples to a category based on which side of the gap they fall on. 12 classes of movements are designed and shown in Fig. 6 (a). The formal 11 classes are meaningful movements while class 12 is meaningless. A multi-class SVM model is trained when given the manually labeled training data. In practice, "one against one" strategy is applied for multi-class SVM and a distribution on all the labels for each movement is the output.

### 3.5. Set-based Soft Discriminative Model

As pointed out that the intra-class variances of our 476 dataset is large since one gesture may have multiple rules 477 to be assembled. In the discriminative model, the rules and 478 479 their corresponding probabilities can be learned from train-480 ing data. Movement sets, both meaningful and meaningless, are considered as rules in our work. Suppose x is the move-481 ments set (observation) and y is the gesture label (hidden 482 state). Normal Discriminative classifiers model the posteri-483 or p(y|x) directly. Suppose each movement in the set has 484 a probability  $P_{w\_mov}$  to be wrongly classified. Therefore, 485

the gesture consisting of movements in such set has a probability of  $P_{w\_ges} = 1 - (1 - P_{w\_mov})^{n_s}$  to be wrongly recognized, where  $n_s$  is the number of movements in one set. Through  $P_{w\_mov}$  may be small,  $P_{w\_ges}$  can be quite large so that it can hardly be tolerated in real-time applications.

In our work,  $P_{w\_mov}$  is relatively small by the proper feature and classifier. Less promotion can be further made. For this reason, a soft discriminative model is proposed to directly decrease the  $P_{w\_ges}$ . We define *m* as the movement with a distribution on all labels, which are denoted as  $M. p_m^M$  represents the probability of movement *m* classified to class *M*. *Ms* represents the trained movement set, also known as assemble rule, which serves as observation.  $G_j$  is the gesture label and serves as hidden state. The probability  $p(G_j)$  is computed as follows:

$$\begin{cases}
Ms_{i} = \langle M_{i1}, M_{i2}, \dots, M_{in_{i}} \rangle i = 1, 2, 3 \dots \\
p(Ms_{i}) = \prod_{k=1}^{n_{i}} p_{m_{k}}^{M_{ik}} \\
p(A_{j}) = \max_{i} \sqrt[n_{i}]{\left[ p(A_{j}|Ms_{i}) \times p(Ms_{i}) \right]}
\end{cases}$$
(2)

where  $n_i$ -th root is used for normalizing since the number of the movements in one set is variant.

Actually, Eq. 2 replaces traditional x with a set Ms and fully use the distribution of classification output  $p_m^M$ . Each rule  $Ms_i$  is evaluated on subsets of a given movement sequence along time. The Ms can be regarded as a sparse joint distribution of the movement in sequence. It provides a soft observation of the discriminative model. The detailed implementation is described by Algorithm 1 and Fig. 7 gives a simple example of this model.

## 4. Experiment

**Dataset and Correct Rate.** We collect a new dataset of interactive presentation gestures that contains five classes: up, down, go, back and click intuitively corresponding to the up, down, left, right and enter in the keyboard. Part of the assemble rules of gestures is shown in Fig. 8, in which movement sets are represented in the form of MHI. The number of movements in one action and the appearance of the same class of MHI are variant, since different subjects act in a quiet different way.

In our dataset, each gesture was performed by three subjects for five times and at two different light conditions: projection light and normal light. Each subject performs three times at a normal speed (4 second/gesture) and the other two times at a fast speed (1 second/ gesture) under each light conditions. The distance between Kinect and subjects is 1.5-2.5m. The depth maps were captured at about 30 frames per second and the size is  $640 \times 480$ . In addition, the PARs are normalized to  $120 \times 260$ . Altogether, the dataset has  $3(subjets) \times 5(times) \times 5(actions) \times 2(illuminations) =$ 

												_
1	Train k	rules	Ms =	$= \{Ms\}$	$s_1, M$	$s_2N$	$Is_k\};$					l
2	repeat			t	1,	2	<i>i</i> 0) /				1	
3	Add a new movement $m_{n+1}$ to a queue of											
	movements;										2	T
4	for	rule i	ndex i	= 1 t	<i>o k</i> <b>d</b>	0						
5		Com	oute th	ie proł	oabili	ty $p(M$	$Is_i$ ) for	all th	e		_	
		sets v	vith eq	ual le	ngth	$n_i$ in q	ueue of				3	
		move	ments	;				1.C \				
6		Multı	ply $p($	$(Ms_i)$	by p	osterio	$\mathbf{r} p(G_j)$	$Ms_i$	;			-
7		Norm	alize	$p(G_j)$	;						Fig	u
8	end	1		$(\alpha$	<b>)</b> .						<b>T</b> 1	1
9	Cho	$\cos e t$	ne max	$\mathbf{x} p(G)$	j); <b>th</b> om						Out	)l( r
10		$(G_j)$	> inre	snoia	unen	hol it i	G				folc	ls
12		Delet	e the r	noven	nents	from (	$G_j,$					
13	else	Defet		noven	ients	nom	lucuc,					
14		No ac	ction is	s detec	cted:							
15	end				,							
16	until no	o mov	ement	is det	ected	any m	ore;					
Ā	Algorith	m 1: (	Gestur	e reco	gniti	on fror	n movei	ment s	sets			
	$\rightarrow$	Mov	vement	ts Sequ	ence	$\rightarrow$	•					
		m1	m2	m3	m4						Tab	ole
	1	0.15	0.12								s n	0
	2	0.13	0.15	p(up)	$=\sqrt[2]{p(}$	up  < 1,2	>) × <i>p</i> (<	1,2 >)			exc	Ц
	3	0.01	0.00	$=\sqrt[2]{0}$	.7 × 0.2	$15 \times 0.15$						
	4	0.12	0.04	= 0.1	255							
	5	0.00	0.16	p(up)	$= \sqrt[2]{p(1)}$	upl < 1.5	$>) \times p(<)$	1.5 > )				
	6	0.10	0.13	$=\sqrt[2]{0}$	$0 \times 0.1$	$5 \times 0.16$	, , , , , , , , , , , , , , , , , , , ,					
	7	0.06	0.09	= 0								
											• 1	
labels Probability distribution 1dat								ti O				
Figure 7. Suppose the priori $p(up  < 1, 2 >) = 0.7$ and $p(up  < 1, 5 >) = 0$ . See from the table, this movement set has the largest joint probability on $< 1, 5 >$ after classification, and has no chance to be labeled as "up". However, in soft discriminative model, this set still has the probability of 0.1255 to be labeled as "up" when $< 1, 2 >$ is chosen.												

150 gestures, 30 samples for each gesture. To test the classification of movement and the recognition of ges-ture, 60 gestures  $(3(subjets) \times 2(times) \times 5(actions) \times$ 2(*illuminations*)) serve as training data and the rest serve as test data. The dataset are labeled manually before train-ing and test. The confusion matrix of 12 classes of move-ments is shown in Table 1. The confusion matrix of 5 class-es of gestures is shown in Table 2. The correct classifica-tion rate of movement achieves 95.23% in 5-folds cross val-



re 8. Samples of five gestures performed by 3 subjects.

e 1. Confusion matrix of 12 classes of movements in the test. method achieves a correct recognition rate of 95.23% in 5-Cross Validation



e 2. Confusion matrix of 5 classes of gestures. None meangesture is those frames. Achieve a correct rate of 90.00% uding "None'



on and the correct recognition rate of gesture achieves 0% in the test.

Complexity reduction. In our work, the gesture deteccomplexity adapts to the arm's movement, that is, the ction frequency is controlled by detection mode derived n PARs and hand position in the previous depth map. In dataset, 18 videos are captured for testing. When coling the dataset, one subject successively performs five different gestures in one video, and each video contains 600 frames. Among these depth frames, only partial of them are selected as active frame for gesture detection. As the results, the computing complexity is vastly reduced, and its reduction ratio is proportional to the number of inactive frame. Fig. 9 (Top) shows the active frame ratio over all frames of 18 videos. The active frame ratio is around 50%, which means that half of depth frames will not calculated for detection. When speaker spend more time on presentation instead of interaction, the active frame ratio will further reduced.

#### CVPR 2013 Submission #1106. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 9. Top: active frames numbers of all the 18 videos. Bottom: interval (blue line) and computational complexity curve of one of our samples. Red line represents our computational complexity while green dash line represents normal computational complexity (Best view in color).

Table 3. Computational complexity in different detection mode. n is the number of pixels in PARs. Chamfer distance maps can be executed in linear (O(n)) time.

Detecting Mode	Processing	Frequency	
Inactive mode	Detecting $(O(n))$	Every 10 frames	
Sami aatiya mada	Detecting $(O(n))$	Every $k$ frames	
Senn-active mode	Detecting $(O(n))$	(1 < k < 10)	
	Detecting $(O(n))$		
A ativa mada	+ Recording $(O(n))$	Every frome	
Active mode	+ Feature	Every frame	
	extraction(O(n))		

Table 3 gives the computational complexity in different detection mode. Fig. 9 (Bottom) shows the complexity from one of the 18 videos, which containing four gestures. At the beginning, hand is detected every 10 frames in an "inactive" mode. Once hand is detected in the PARs and is close enough to the screen, the detecting mode is switched to "active". If the hand leaves the PARs or become far away from screen, the number decreases gradually. See from Table 4, the active frame ratio reduces to 46.47% and the correct recognition rate does not decrease much. This method can be potentially used for wireless transfer of frames.

697 Comparison on Features. In the field of gesture recognition, the trajectories of gestures are always represented as
699 a set of points (e.g., sampled positions of the head, hand, and eyes) in a 2-D space before being decomposed. For
701 example, HoGS is a descriptor proposed by Sadeghipour

Tal	ble	4.′	The to	otal fram	e num	ber an	d corr	ect recog	gnition	rate.	(AF-	
_			_	_		~ ~ ~	-	_		_		



Figure 10. (a) Comparison on two features. (b) Trajectory has problem with outliers (clothes), the straight line is wrongly connected. (c) The five attributes for trajectory in our experiment.

et al. [11]. They combine this feature with SVM to solve the challenging problem of gesture-based object recognition. Though trajectory is obtained by tracking hand in the first place, MHI is our final feature. The reason lies in two folds. First, compared with sensitive point detection, the method to generate MHI is more robust since it simulates original silhouettes. Second, the MHI implicitly represents the direction of movement while a trajectory of hand points has less information of that information. To compare the two features, an experiment using trajectory as feature is conducted on our dataset.

As trajectories have been segmented by MLS, some attributes can be extracted from the curve of segmentations like the method in [11]. Five attributes are used: height, width, length, orientation and center of the curve (see Fig. 10 (c)). Combined with SVM, this feature has a low movement correct rate and gesture correct rate (see Table 5). As stated above, the main reason is the sensitiveness of points on trajectories. For example, pictures (a) and (b) in Fig. 4 are the comparisons between trajectory feature and MHI feature. In (a), the two features have almost the same discriminating power and are both correctly classified in experiment. In (b), since the clothes of the subject used to enter the bounding box and produce some outliers at the left bottom region, trajectory fails to describe this movement because of a wrongly connected straight line while the MHI feature is correctly classified. That mainly owns to the abundant original information the MHI feature contains.

Set-based Soft Discriminative Model. On the training data set including 60 samples, a posterior p(G|Ms) is trained. In traditional discriminative model, Ms includes the movement sets that occur at least once in the training

762

763

764

765

766

767

768

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

Table 5. Comparison result. Movement correct classification rate
are computed by 5-folds CV. (DM: Discriminative Model; MCCR:
Movement Correct Classification Rate; GCRR: Gesture Correct
Recognition Rate ).

Feature	Model	MCCR	GCRR
MHI + Semi-circle	Traditional DM	95.23%	76.67%
MHI + Rectangle	Soft DM	79.31%	76.67%
Trajectory	Soft DM	48.32%	62.23%
MHI + Semi-circle	Soft DM	95.23%	90.00%

data. New movement sets have no chance to be changed to 769 the nearest one in the Ms while soft discriminative model 770 does. That is because soft discriminative model and choose 771 the best one according to the distribution of movements  $m_i$ 772 (see Eq. 2). The traditional model is also tested on the test 773 set including the rest 90 samples and the correct recognition 774 775 rate is 76.67%. This result shows that the correct recognition rate is limited by the scale of train set since traditional 776 discriminative models directly model a conditional distri-777 bution over gesture labels given the observations. Besides, 778 our observation is set-based, some observation sets on test 779 set never occurs on training set. In our approach, this obser-780 vation set is mapped to the one exists on the soft discrim-781 inative model trained by small scale train set. After using 782 the soft discriminative model, the correct recognition rate is 783 90.00% (see Table 5). 784

## 5. Conclusions and Future Work

We propose a framework for gesture recognition and test it in our dataset. Experimental results show that our method is efficiency due to the detection mode transfer and robust due to the MHI feature and soft discriminative model. In addition, depth maps can further compressed in our framework, which will enhance the efficient. To recognize gesture from compressed images is our future work.

# References

- M. Blank, L. Gorelick, E. Shechtman, M. Irani, and R. Basri. Actions as space-time shapes. In *ICCV*, volume 2, pages 1395–1402. IEEE, 2005. 2
- [2] A. Bobick. Movement, activity and action: the role of knowledge in the perception of motion. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 352(1358):1257– 1265, 1997. 3
- [3] A. Bobick and J. Davis. The recognition of human movement using temporal templates. *PAMI*, 23(3):257–267, 2001. 4
- [4] X. Feng and P. Perona. Human action recognition by sequence of
  movelet codewords. In *3D Data Processing Visualization and Trans- mission*, pages 717–721. IEEE, 2002. 3

- [5] K. Fujimura and X. Liu. Sign recognition using depth image streams. In FG, pages 381–386. IEEE, 2006. 2
- Y. Jia, C. Huang, and T. Darrell. Beyond spatial pyramids: Receptive field learning for pooled image features. In *CVPR*, pages 3370–3377.
   IEEE, 2012. 5
- [7] A. Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. *Advances in neural information processing systems*, 14:841, 2002. 3
- [8] A. Klaser, M. Marszalek, and C. Schmid. A Spatio-Temporal Descriptor Based on 3D-Gradients. In M. Everingham, C. Needham, and R. Fraile, editors, *BMVC*, pages 275:1–10, Leeds, Royaume-Uni, 2008. British Machine Vision Association. 3
- [9] I. Laptev. On space-time interest points. *IJCV*, 64(2):107–123, 2005.
   3
- [10] R. Poppe. A survey on vision-based human action recognition. *IVC*, 28(6):976–990, 2010.
- [11] A. Sadeghipour, L. Morency, and S. Kopp. Gesture-based object recognition using histograms of guiding strokes. *BMVC*, 2012. 2, 3, 7
- [12] C. Schuldt, I. Laptev, and B. Caputo. Recognizing human actions: A local svm approach. In *ICPR*, volume 3, pages 32–36. IEEE, 2004.
   2
- [13] P. Scovanner, S. Ali, and M. Shah. A 3-dimensional sift descriptor and its application to action recognition. In *Proceedings of the 15th international conference on Multimedia*, pages 357–360. ACM, 2007. 3
- [14] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *CVPR*, pages 1297–1304. IEEE, 2011. 2
- [15] Y. Song, D. Demirdjian, and R. Davis. Tracking body and hands for gesture recognition: Natops aircraft handling signals database. In *FG*, pages 500–506. IEEE, 2011. 2, 3
- P. Turaga, R. Chellappa, V. Subrahmanian, and O. Udrea. Machine recognition of human activities: A survey. *Circuits and Systems for Video Technology, IEEE Transactions on*, 18(11):1473–1488, 2008.
- [17] J. Wang, Z. Liu, Y. Wu, and J. Yuan. Mining actionlet ensemble for action recognition with depth cameras. In *CVPR*, pages 1290–1297. IEEE, 2012. 2
- [18] Q. Wang, X. Chen, and W. Gao. Skin color weighted disparity competition for hand segmentation from stereo camera. In *BMVC*, pages 66–1, 2010. 2
- [19] D. Weinland, E. Boyer, and R. Ronfard. Action recognition from arbitrary views using 3d exemplars. In *ICCV*, pages 1–7. IEEE, 2007.
   3
- [20] G. Willems, T. Tuytelaars, and L. Van Gool. An efficient dense and scale-invariant spatio-temporal interest point detector. *ECCV*, pages 650–663, 2008. 3