

Distributed Spatio-Temporal Gesture Recognition in Sensor Arrays

HOMA HOSSEINMARDI, University of Colorado Boulder
AKSHAY MYSORE, University of Colorado Boulder
NICHOLAS FARROW, University of Colorado Boulder
NIKOLAUS CORRELL, University of Colorado Boulder
RICHARD HAN, University of Colorado Boulder

We present algorithms for gesture recognition using in-network processing in distributed sensor arrays embedded within systems such as tactile input devices, sensing skins for robotic applications, and smart walls. We describe three distributed gesture recognition algorithms that are designed to function on sensor arrays with minimal computational power, limited memory, limited bandwidth, and possibly unreliable communication. These constraints cause storage of gesture templates within the system and distributed consensus algorithms for recognizing gestures to be difficult. Building up on a chain vector encoding algorithm commonly used for gesture recognition on a central computer, we approach this problem by dividing the gesture dataset between nodes such that each node has access to the complete dataset via its neighbors. Nodes share gesture information among each other and then each node tries to identify the gesture. In order to distribute the computational load among all nodes, we also investigate an alternative algorithm where each node that detects a motion will apply a recognition algorithm to part of the input gesture and then will share its data with all other motion nodes. We then show that a hybrid algorithm that distributes both computation and template storage can address trade-offs between memory and computational efficiency.

Categories and Subject Descriptors: C.2.2 [Computer-Communication Networks]: Network Protocols

General Terms: Design, Algorithms, Performance, Reliability

Additional Key Words and Phrases: distributed computing; error resilient; gesture recognition

ACM Reference Format:

Homa Hosseinmard, Akshay Mysore, Nicholas Farrow, Richard Han, Nikolaus Correll, 2015. Distributed Spatio-Temporal Gesture Recognition in Sensor Arrays. *ACM Trans. Autonom. Adapt. Syst.* V, N, Article A (March 2015), 20 pages.

DOI : <http://dx.doi.org/10.1145/2744203>

1. INTRODUCTION

Gesture-based control, ranging from simple directional swipes to input of complex characters has emerged as a standard method for human-computer interaction. Most systems however are spatially limited to sensor arrays, e.g., tactile displays, range-finding cameras such as the Xbox Kinect, and require information processing in a central unit. We wish to make gesture recognition available to a wider range of surfaces, which embed memory, computation and communication capabilities and can therefore function without a central processing unit. Examples of such surfaces range from a building wall made of “smart bricks”, robotic sensing skins, or “smart paint” [Butera

This work has been supported by NSF awards #1153158 and #1150223. We are grateful for this support.

Author's addresses: H. Hosseinmardi, A. Mysore, N. Farrow, N. Correll, R. Han, Computer Science Department, University of Colorado Boulder, 430 UCB, Boulder, CO 80309

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2015 ACM 1556-4665/2015/03-ARTA \$15.00

DOI : <http://dx.doi.org/10.1145/2744203>

2002; McEvoy and Correll 2015]. In the short term, we are interested in smart building materials that allow users to control environmental conditions using gestural input without a central computer processing the data and which can do so independently of the arrangement of the smart bricks. An example of simple swipe-based interaction with such a system is shown in Figure 1, which allows a user to use linear or circular gestures to control specific functions. In this previous work, a user interacts with specific hardware functions through submenus providing slider-like or dial-like input for room temperature or brightness control, e.g. [Farrow et al. 2014]. While such an approach provides a wide range of possible interactions by hard-coding distributed controllers to recognize up-down and circular swipes, limiting gestures to simple swipes quickly becomes non-intuitive as more functions become available, for example if this interface should be used to control multiple aspects of a building, such as lighting, heating, ventilation and cooling, or even selecting amongst different ambient music channels [Farrow et al. 2014]. Adding the capability to recognize actual letters and other complex gestures should therefore drastically increase the vocabulary and extensibility of such a system.

In the long run, we envision a new class of smart materials that completely embeds and integrates sensing, computation, and actuation [Paradiso et al. 2004]. Such a hardware base could use the presented distributed pattern recognition algorithms to allow for rich interaction with their users. For example, when embedded in clothes [Nakad et al. 2003; Profita et al. 2015], gestures might allow communication with other devices such as music players or cell phones in an intuitive manner. Integration into a robotic skin [Lumelsky et al. 2001; Hughes and Correll 2014] might allow users to use tactile gestures to control the robot. Finally, materials with shape changing capabilities [McEvoy et al. 2013; McEvoy and Correll 2014] could enable furniture that can adapt to their user's needs via rich gestural input.

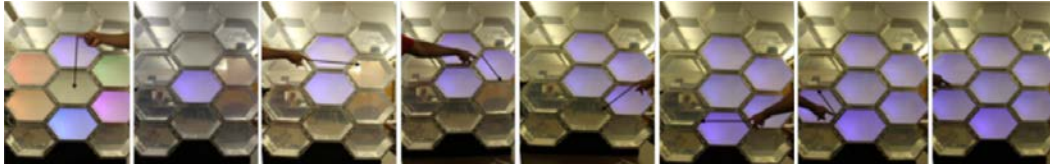


Fig. 1: An example simple interaction with a Smart façade in our prior work: touching any of the wall element will provide a color coded “menu” whose items are selected by dragging them into the center. In this example, selecting the blue menu item switches to a dial-like submenu where input is then provided via a circular gesture [Farrow et al. 2014].

In previous work, we have developed multi-cast communication algorithms to limit communication in a distributed system to devices that are actually part of the gesture [Hosseinmardi et al. 2012; Ma et al. 2012], though that previous work did not investigate gesture recognition algorithms. This paper is a first attempt to study distributed gesture recognition with a view towards supporting next-generation smart materials. A key challenge here is that all nodes that are part of the gesture need to exchange information with all other nodes, which is difficult in harsh communication environments with minimal bandwidth. Emphasis on minimalist computation, memory, and communication might enable distributed gesture recognition on other surfaces where sensors and computation are tightly integrated. In such systems, links are likely to

eventually fail due to imperfect low cost fabrication and wear and tear of the materials, thus providing additional complications to a distributed gesture recognition algorithm.

In this paper, we introduce three distributed algorithms that are adapted to minimize memory, minimize computation, and a hybrid algorithm which minimizes both memory and computation. We show that the hybrid algorithm achieves performance that is comparable to the distributed computation algorithm in terms of recognition reliability and robustness to packet-loss, but maintains the memory efficiency of the distributed dataset algorithm.

2. RELATED WORK

There exists a large body of work on gesture recognition and online handwriting recognition in touchpads, tablets, tablet computers [Wu and Huang 1999; Mitra and Acharya 2007] or other electronic devices [Davis and Lyall 1986]. Although the gestures we are interested in are not necessarily limited to alphabet characters, stroke-based character recognition is a mature field and allows encoding and recognizing a large variety of gestures. These algorithms deal exclusively with a centralized representation of the gesture, that is gestures are recorded by a specific device such as a camera, range-finding device or inertial-measurement unit (IMU) [Schlömer et al. 2008], co-located with a computer and processed thereon. In contrast, our paper seeks to implement recognition of complex gestures in a fully distributed and scalable fashion.

[Nagi et al. 2012] is most closely related in terms of distributing the gesture recognition problem and introduces a gesture recognition algorithm to overcome the perception limitations of simple swarm robots. Robots observe a simple hand gesture from different angles, attempt its classification, and then share their posterior probability vectors and arrive via a distributed consensus at a decision. This approach does not consider distributing memory or communication. Rather, it capitalizes on the variations in perception that each robot incurs by looking at the gesture from a different angle. In our work instead, we aim to distribute computation and memory in order to perform classification, while being robust to communication failures which may make individual units perceive a gesture slightly differently.

A common approach to encoding gestures/hand-written characters is to use a “chain code” taking into account its trajectory [Elmezain et al. 2009; Ozer et al. 2001; Confer and Chapman 2004]. In fact, [Elmezain et al. 2009] shows that features that convey trajectory information of the swipe are the most robust features for gesture/character recognition. For this, consecutive points of a gesture are sampled at regular intervals, and a discrete number, e.g., 0-7 to encode angles in 45 degree intervals, are assigned based on the angle of the two consecutive samples. The resulting number streams can then be classified, e.g., using Hidden Markov models (HMM) [Lee and Kim 1999; Wilson and Bobick 1999; Elmezain et al. 2009], Dynamic Time Warping (DTW) [Sakoe and Chiba 1990; Carmona and Climent 2012], support vector machines [Bahlmann et al. 2002] or neural networks [Murakami and Taguchi 1991].

While these approaches are intellectually appealing, they suffer from high computational complexity that make them less suitable for resource-constrained distributed sensor arrays in smart materials. An alternative is to employ a simple and robust recognition method that simply compares the Euclidian distance of a sequentially acquired chain code [Confer and Chapman 2004]. Although this approach requires the users to provide the gestures in the correct direction, it has been successfully used on platforms such as the Palm Pilot user interface. This approach further exhibits the advantage that it can be modified to allow us to divide both the dataset and computational load across a grid of distributed nodes. We selected this algorithm over a HMM

or DTW approach due to its computational simplicity and its proven robustness, which makes it potentially applicable for implementation in distributed architectures with very limited computational capabilities. For example, DTW has computational complexity $O(n^2)$ since it uses a dynamic programming algorithm to compute the distance between two chain codes, while [Confer and Chapman 2004] use simply Euclidean distance with $O(n)$ to compute the distance, where n is the dimension of feature vector. We also note that distributing DTW is not straightforward as the underlying dynamic programming approach requires all parts of the dataset to provide an optimal estimate.

The idea to create user interfaces in materials that are equipped with distributed computation has been first articulated by [Lifton et al. 2002] who introduced “Push-pin computing”, and the idea of a smart, distributed flooring system [Schlömer et al. 2008] as an architecture for algorithms such as proposed in this paper. Yet, these ideas have become economically feasible only recently with microcontrollers becoming more powerful, smarter and cheaper, which we believe will lead to a renewal of interest in distributed algorithms that can operate on such platforms.

3. BASIC GESTURE RECOGNITION ALGORITHM

We consider a distributed array of sensing nodes arranged as a lattice, with each node in the lattice having limited computational abilities and local connectivity to its immediate neighbors. We further assume that all nodes involved in a gestural event can communicate with each other by multi-hop multi-cast communication [Hosseinmardi et al. 2012; Ma et al. 2012]. In this manner each node involved in the gesture, in short a *motion node*, will create a table of other group members’ unique identification number *ID*, and their relative coordinates and temporal order that can be inferred from the network topology, e.g., a square grid or a hexagonal arrangement.

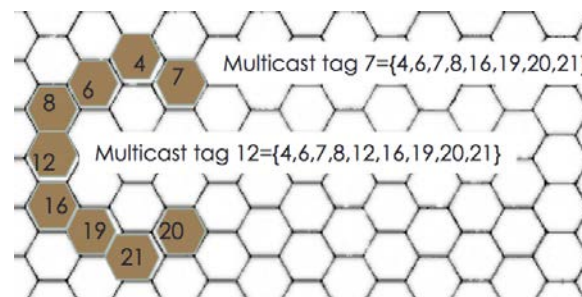


Fig. 2: The shaded nodes have detected the motion gesture with the order $\{7, 4, 6, 8, 12, 16, 19, 21, 20\}$. They will start sharing their packet ID. All motion nodes and their immediate neighbors will help in forwarding data packets of other motion nodes, however some packets may not be delivered correctly (data packet of node 12 has not been delivered to node 7). Finally, node 7 has a table including information from nodes $\{4, 6, 7, 8, 16, 19, 20, 21\}$, and node 12 has a table including information from nodes $\{4, 6, 7, 8, 12, 16, 19, 20, 21\}$

As in centralized gesture recognition, pattern recognition will be based on comparison of an input stream with a dataset of template vectors. A key challenge here is that gestures likely have different sizes and therefore require scaling as a pre-processing

step. In our scenario, gesture size is determined by the number of motion nodes detecting such a gesture. Therefore, data vectors inside motion node tables need to be interpolated or extrapolated to match the template length. Classification can then be performed by finding the template that minimizes a distance metric with the measurement vector. While there are multiple ways to encode the gesture and perform training and classification, the algorithms described in this paper build up on [Confer and Chapman 2004], who proposed a centralized algorithm that uses the Euclidian distance between sequential chain codes. We chose this algorithm due to its linear computational complexity, proven robustness on a full-size alphabet, and relative ease to distribute both its memory and computational requirements.

In order to create a feature vector in a distributed system, a motion node emits a data packet upon sensing an event such as touch or light. Such an event packet conveys both the spatial position information of the node that has detected an event as well as temporal information about the time the event was recorded. This information is then used on each node receiving the packet to reconstruct a complete spatio-temporal feature vector from all the collected packets.

After describing pre-processing steps that are common to all algorithms presented in this paper, we propose methods for distributing memory, distributing computation, and combining both into a hybrid method.

3.1. Feature Extraction

As shown in Figure 2, when someone draws a gesture, each motion node that senses the gesture also needs to learn about other nodes that sense the motion. As we describe in [Hosseinmardi et al. 2012; Ma et al. 2012] after detecting an event, each motion node will share its ID , spatial and order information with other motion nodes. Each node gathers the spatio-temporal information into a table that will be used later to extract feature vectors.

Because different directions of input gestures might have different meanings, all nodes need to know the order in which events were recorded. Nodes will sort member vectors and their spatial information based on the gesture temporal order for the purpose of feature extraction. Each node m will retrieve a certain number of gesture points from collecting shared packets by other nodes sensing the same gesture. The node will maintain three data vectors as motion packets are collected: ID^m , g_x^m and g_y^m where the superscript indicates that these vectors are collected by node m . The g_x^m and g_y^m vectors correspond to the spatial coordinates of the sensors that detected the motion. In order to compare gesture vectors to a dataset of possible gesture templates, a new sequence of vectors will need to be computed into a vector chain whose set size N is of equal size to the templates in the dataset. For this purpose each node m will interpolate g_x^m and g_y^m vectors to obtain x^m and y^m vectors which now describe a unique feature gesture composed of N chain vectors connecting $N + 1$ ordered points. For example, in Figure 3, any number of gestural sensing events may be collected at a particular node; these are interpolated to produce an $N + 1$ point gesture.

Considering that users may draw gestures with different sizes thereby spanning different numbers of nodes, and that packet loss may cause different group members to have different numbers of nodes in their table, input gestures are likely to have different vector lengths. The number of motion node packets received by m will be denoted with M . Since some of the packets might not be properly recorded or delivered, each node has a potentially different number of gesture points M , and consequently a different view of the same gesture. However, we employ the chain technique to extract a fixed-length feature vector on each node independantly. In this way, the fixed-sampling technique is fairly robust within limits both to how small a gesture is made, as well

as to packet loss or sensor failure. Figure 4 shows three chain vectors with $N = 10, 16,$ and 32 .

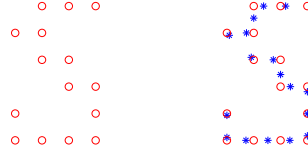


Fig. 3: Input curve of arbitrary size at a motion node, where red circles are sensors coordinates (left) and asterisks are used to show the interpolated points to get a fixed gesture length $N + 1$ (right).

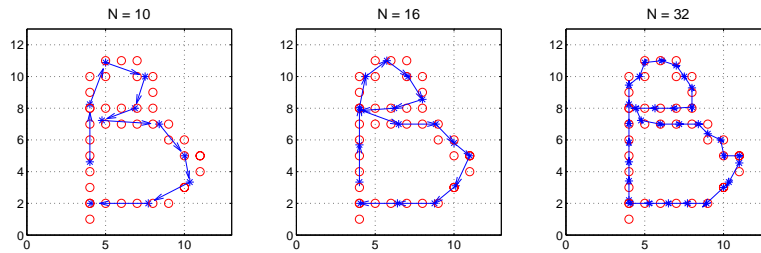


Fig. 4: Chain vector for different N in a 12×12 network

A gesture can now be represented by the N vectors that are all the same length connecting an interpolation of the sample points. We denote a chain of vectors recorded by node m as $V_m = \{\vec{v}_1^m \dots \vec{v}_N^m\}$

$$\vec{v}_i^m = \ell_x^m(i)(1 \ 0)^T + \ell_y^m(i)(0 \ 1)^T \quad (1)$$

Here, each element of $\ell_x^m(i)$ and $\ell_y^m(i)$ is calculated from the (x_i^m, y_i^m) and (x_{i+1}^m, y_{i+1}^m) coordinates of the points after interpolation

$$\ell_x^m(i) = x_{i+1}^m - x_i^m, \quad \text{where } i = 1, 2, \dots, N \quad (2)$$

$$\ell_y^m(i) = y_{i+1}^m - y_i^m, \quad \text{where } i = 1, 2, \dots, N \quad (3)$$

where i and $i + 1$ are adjacent points and

$$\mathbf{x}^m = (x_1^m, x_2^m, \dots, x_{N+1}^m) \quad (4)$$

$$\mathbf{y}^m = (y_1^m, y_2^m, \dots, y_{N+1}^m) \quad (5)$$

3.2. Basic Matching Algorithm

In order to recognize what gesture is made, we use a 1-nearest neighbor (1-NN) clustering algorithm to find the template s that is closest to the gesture m .

Let $\delta_{ms} = |V_m - V_s|$ be the distance between a vector chain on node m and template $s \in \mathcal{DS}$, a dataset containing all templates. The function $f(V_m)$ will return the lowest distance between input gesture vector chain on node m and dataset template chains s , such that

$$f(V_m) = \min_s |V_m - V_s| \quad (6)$$

Here, $|V_m - V_s| = \sum \|\vec{v}_i^m - \vec{v}_i^s\|$ is the distance between two chain vectors and $\|\cdot\|$ is the Euclidean distance between two vectors. Algorithm 1 summarizes our basic matching algorithm.

ALGORITHM 1: Each node m will compute the total chain vector difference between its input vectors and all template vectors.

- 1: Wait until the gesture is complete
 - 2: **if** it is a motion node m **then**
 - 3: Sort table
 - 4: Interpolate sensor coordinates to find $N + 1$ points along the gesture
 - 5: Create $V_m = \{\vec{v}_1^m \dots \vec{v}_N^m\}$ using equations 1-3
 - 6: $\delta_{ms} = |V_m - V_s|$, for all $s \in \mathcal{DS}$
 - 7: $s^* = \arg \min_s \delta_{ms}$ and $s = 1, \dots, |\mathcal{DS}|$
 - 8: $\tilde{c} = \text{character of } (s^*)$
 - 9: goto 1
 - 10: **end if**
-

4. DISTRIBUTED GESTURE RECOGNITION ALGORITHMS

In this section we propose three algorithms that achieve distributed gesture recognition: 1) Distributed Dataset (DD), 2) Distributed Vector computation (DV) and 3) Hybrid Approach (Distributed Dataset and Vector computation, DDV).

4.1. Distributed Dataset Algorithm

In order to match a gesture to a known pattern, all possible patterns need to be stored somewhere in the system. The dataset \mathcal{DS} includes a specific number of templates V_s for each possible input. The size $|\mathcal{DS}|$ depends on the granularity of the pattern and the number of possible patterns that can be recognized. Dividing \mathcal{DS} into K subsets can therefore lead to considerable memory savings and the only solution to make high granularity / high variability pattern recognition feasible on a distributed system with minimal memory. In this case, memory requirements are traded off with additional communication requirements for nodes to retrieve dataset information.

In order to minimize this communication overhead, we propose storing the partial datasets ds_n in node n in such a way that every node m has access to the full dataset via its one-hop neighborhood $\mathcal{N}(m)$, i.e.

$$\mathcal{DS} = \bigcup_n \{ds_n | n \in \mathcal{N}(m) \cup m\} \quad (7)$$

We can now compute the distances between the vectors of the input gesture with the templates contained in the partial dataset ds_n . In this manner neighbors $\mathcal{N}(m)$ of each motion node m will help in computing the distances from other parts of a

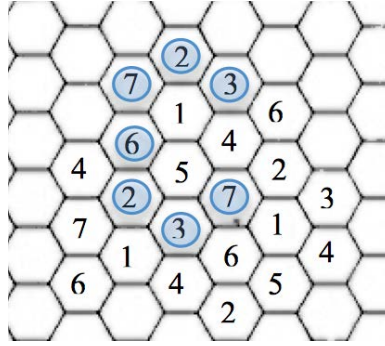


Fig. 5: Possible distribution of pattern datasets between nodes. Note that the dataset has been split into $K = 7$ subsets, with each node storing only one subset. Motion nodes corresponding to the character “C” are highlighted.

dataset. Since the template s that has the minimum distance from the input vector chain will be chosen as the output \hat{g} , each neighbor node $n \in \mathcal{N}(m)$ will send just the minimum distance and the corresponding template number s ($s, \delta_n = \min_s \delta_{ns}$) to its group members.

Consider Figure 5 where a dataset has been distributed in $K=7$ parts and all nodes have access to the complete dataset via their one-hop neighbors. For an input “C” drawn on the network, we can see that the gesture has covered nodes with ds_2, ds_3, ds_6 and ds_7 , but did not extend to nodes containing ds_1, ds_4 and ds_5 . For each node m , its neighbor nodes that have different parts of the gestural dataset assist in calculating the distance from m 's vector V_n to their templates V_s and share that information back with node m . At this point, m can choose the template that has the smallest δ_n .

Since all neighbor groups contain the complete dataset, instead of sharing δ_n among group members, it will be shared between one-hop neighbors. In this way, all neighbor nodes $n \in \mathcal{N}(m)$ will send their (s, δ_n) to node m . Node m will compare all the distances δ_k for $k \in K$, find the minimum of them, and select template number s corresponding to the minimum δ_n as the output \tilde{c} , i.e. the closest gesture in the dataset, which is the consensus computational decision of the distributed set of nodes. The pseudo code for the distributed dataset (DD) is provided in Algorithm 2.

As we will later show, distributing the dataset will result in dramatically lower robustness for recognizing patterns under the influence of packet loss. The reason is that when neighbor nodes send δ_k to node m , if the packet that includes the best template match is not delivered to m , it will impact heavily on the decision of node m .

4.2. Distributed Vector computation

Distributing the dataset already distributes computation somewhat — distances are only calculated for the local part of the dataset ds_n — but we are also interested in further distributing computation to reduce computational requirements on any single node in the sensor array. The basic idea is that each node computes distances only for a subset of the vectors that make up a full gesture. The results of these computations can then be shared with their neighbors.

First, each node will compute the pieces of data with length $\lceil N/M \rceil$ in the interval $[a, b]$. The boundaries of this interval are calculated as follows for the p^{th} node:

ALGORITHM 2: Distributed Dataset (DD): Each node n has access to a subset of dataset ds_n . The dataset is distributed between nodes such that for each motion node m

$$\mathcal{DS} = \bigcup_n \{ds_n | n \in \mathcal{N}(m) \cup m\}$$

```

1: Wait until the gesture is complete
2: if it is a neighbor node  $n$  then
3:   Sort its table, find motion nodes which are on its one-hop neighborhood
4:   Interpolate sensor coordinates to find  $N + 1$  points along the gesture
5:   for all motion node  $m$  in the neighborhood of node  $n$  do
6:     Create  $V_m = \{\vec{v}_1^m \dots \vec{v}_N^m\}$  using equation 1-3
7:      $\delta_{n,s} = |V_n - V_s|$  for all  $s \in ds_n$ 
8:      $\delta_n = \min_s \delta_{n,s}$ ,  $s = \arg \min_s \delta_{n,s}$ 
9:      $P_n \leftarrow (\text{character of } (s), \delta_n)$ 
10:    Send packet  $P_n$  to neighbor motion node  $m$ 
11:  end for
12:  go to step 1
13: end if
14: if it is a motion node  $m$  then
15:   Sort its table
16:   Interpolate sensor coordinates to find  $N + 1$  points along the gesture
17:   Create  $V_m = \{\vec{v}_1^m \dots \vec{v}_N^m\}$  using equation 1-3
18:    $\delta_{m,s} = |V_m - V_s|$  for all  $s \in ds_m$ 
19:   Receive packets from its neighbors
20:    $n^* = \arg \min_n \delta_n$ , where  $n \in \mathcal{N}(m) \cup m$ 
21:    $\tilde{c} = \text{character in the packets from node } n$ 
22:   go to step 1
23: end if

```

$$a = (p - 1) \times \lceil N/M \rceil + 1 \quad (8)$$

$$b = (p - 1) \times \lceil N/M \rceil + \lceil N/M \rceil = p \times \lceil N/M \rceil \quad (9)$$

$$V_m^p = \{\vec{v}_i^m | i \in [a, b]\} \quad (10)$$

and

$$V_s^p = \{\vec{v}_i^s | i \in [a, b]\} \quad (11)$$

where $\lceil \cdot \rceil$ results in the first integer number equal or bigger than its argument. For motion node m , M is the size of its group. Since the size N/M is not always an integer number, the last motion node may have vector size less than $\lceil N/M \rceil$. In this way vector V_m is divided to $N_{cmpt} = \frac{M}{\lceil N/M \rceil}$ and $p = 1, \dots, N_{cmpt}$. Figure 6 illustrates a possible segmentation of N elements of a gestural chain vector into M subsets.

Each node will calculate the distance of its partial vector V_m^p from dataset templates partial vector V_s^p , $\delta_{m,s}^p = |V_m^p - V_s^p|$ and send it out to all other group members. Each node m will put all $\delta_{m,s}^p$ together for each template s to construct a complete distance vector $\delta_{m,s} = [\delta_{1,s}^p \dots \delta_{M,s}^p]$.

While providing a linear decrease in computational time per node, this method requires more communication packets in order to share partial computation results, which would seem to make the algorithm more susceptible to packet loss. However, as we will show, the distributed vector (DV) algorithm is actually more resilient to packet loss than the distributed dataset approach, because packet loss in the distributed com-

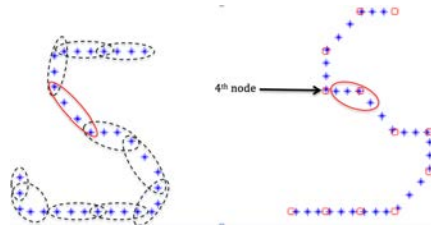


Fig. 6: Possible segmentation of a shapes into elements of length $\lfloor N/M \rfloor$. Each node will calculate the distance of the partial vector of its observation (right) from the partial vector of dataset templates (left).

putation approach only jeopardizes information about parts of a gesture. For example, if a particular motion node receives all but one of the partial comparisons from the other motion nodes, it will still have a fairly good estimate of the closest template gesture to the sensor data, since it will have all comparison segments but one. Pseudocode for (DV) computation is provided in Algorithm 3.

ALGORITHM 3: Distributed Vector (DV): Each node n will compute the distance of part of its input vector from corresponding part of all template vectors.

- 1: Wait until the gesture is complete
 - 2: **if** it is a motion node m **then**
 - 3: Sort its table, find its order p
 - 4: Interpolate sensor coordinates to find $N + 1$ points along the gesture
 - 5: $a = (p - 1) \times \lfloor N/M \rfloor + 1$
 - 6: $b = p \times \lfloor N/M \rfloor$
 - 7: $V_m^p = \{\vec{v}_i^m \mid i \in [a, b]\}$
 - 8: $\delta_{ms}^p = |V_m^p - V_s^p|$ for all $s \in \mathcal{DS}$
 - 9: $P_m \leftarrow (m, [\delta_{m1}^p \dots \delta_{m|\mathcal{DS}|}^p])$
 - 10: Share packet P_m with other motion nodes
 - 11: $\delta_{ms} = \sum_{p=1}^{N_{cmpt}} \delta_{ms}^p$ for all $s \in \mathcal{DS}$
 - 12: $s^* = \arg \min_s \delta_{ms}$ for all $s \in \mathcal{DS}$
 - 13: $\tilde{c} = \text{character of } (s^*)$
 - 14: go to step 1
 - 15: **end if**
-

4.3. Hybrid Method

Combining distributed data storage and distributed computation allows a trade off whereby the advantages of one technique mitigate the drawbacks of the other method. The distributed dataset approach is memory efficient but lacks robustness against packet loss, while the distributed computational vector approach is more computationally efficient and more robust to packet loss, but has large memory requirements when the dataset is not distributed.

We propose a hybrid distributed data and vector computation (DDV) algorithm that combines both approaches. Each motion node will compute the distance of V_n^p from part of the dataset templates ds_n . Neighbors $n \in \mathcal{N}(m)$ will send the distances of V_n^p from other parts of the dataset to the motion node m . Then motion node m will share

the distance of V_n^p from the whole dataset with all other group members. Pseudo-code for the hybrid DDV algorithm is shown in Algorithm 4.

ALGORITHM 4: Distributed Dataset and Vector (DDV): Each node has m access to a subset of dataset ds_m and will compute the distance of part of its input vector to the corresponding part of all template vectors for of ds_m . The dataset is distributed between nodes such that for each motion node m

$$\mathcal{DS} = \bigcup_n \{ds_n | n \in \mathcal{N}(m) \cup m\}$$

```

1: Wait until the gesture is complete
2: if it is a neighbor node  $n$  then
3:   Sort its table, find motion nodes which are in its one-hop neighborhood
4:   Interpolate sensor coordinates to find  $N + 1$  points along the gesture
5:   for all motion node  $m$  with order  $p$  in the neighborhood of node  $n$  do
6:      $a = (p - 1) \times \lceil N/M \rceil + 1$ 
7:      $b = p \times \lceil N/M \rceil$ 
8:     Create partial vector of vectors:  $V_n^p = \{\vec{v}_i^n | i \in [a, b]\}$ 
9:      $\delta_{n,s}^p = |V_n^p - V_s^p|$  for all  $s \in ds_k$ 
10:     $P_n \leftarrow (m, k, [\delta_{n,1}^p \dots \delta_{n|ds_k|}^p])$ 
11:    Send packet  $P_n$  to neighbor motion node  $m$ 
12:  end for
13:  go to step 1
14: end if
15: if it is a motion node  $m$  then
16:   Sort its table, find its order  $p$ 
17:   Interpolate sensor coordinates to find  $N + 1$  points along the gesture
18:    $a = (p - 1) \lceil N/M \rceil + 1$ 
19:    $b = p \times \lceil N/M \rceil$ 
20:    $V_m^p = \{\vec{v}_i^m | i \in [a, b]\}$ 
21:    $\delta_{m,s}^p = |V_m^p - V_s^p|$  for all  $s \in ds_n$ 
22:   Receive packets from its neighbor and construct  $[\delta_{m,1}^p \dots \delta_{m|DS|}^p]$ 
23:    $P_m \leftarrow (m, [\delta_{m,1}^p \dots \delta_{m|DS|}^p])$ 
24:   Share packet  $P_m$  with other motion nodes
25:    $\delta_{m,s} = \sum_{p=1}^{N_{cmpt}} \delta_{m,s}^p$  for all  $s \in \mathcal{DS}$ 
26:    $s^* = \arg \min_s \delta_{m,s}$  for all  $s \in \mathcal{DS}$ 
27:    $\tilde{c} = \text{character of } (s^*)$ 
28:   go to step 1
29: end if

```

As a comparison of the three different techniques, we summarize in Table I the order of savings achieved in computation and space of DD, DV and DDV for a typical motion node with a group of size M , K nodes in a neighborhood, and a size of $|\mathcal{DS}|$ for the data set.

5. EVALUATION

We evaluate the accuracy and performance of each of the three distributed algorithms in a Matlab simulation. We used a 3-fold cross validation technique to evaluate the algorithms. We created 21 templates for each of the English alphabet's 26 characters. We will first evaluate the distributed algorithms and compare them to the Basic algorithm in a typical scenario of a 10x10 network and vector length $N = 16$. As we will see, the Hybrid DDV algorithm outperforms both the DD and DV algorithms. Therefore, we

Table I: Order of savings in computation and space (memory) for motion nodes assuming neighborhoods of size K and group sizes M (N_{cmpt} was defined in section 4.2).

	DD	DV	DDV
space	K	1	K
computation	K	N_{cmpt}	$K N_{cmpt}$

select it as the one for the remainder of the comparisons and evaluate the effect of other parameters such as network size and feature vector length on the performance of the Hybrid DDV algorithm. Finally, we consider applications that require smaller alphabets and how we can choose the subset of gestures that maximize accuracy.

5.1. Data Collection

In order to gather the dataset we used a touch screen iPad as the input interface. We collect 21 different templates (size and shape) of gestures for each of the 26 characters for a total of 546 templates in the dataset, Figure 7 left. We have trained examples provided by a set of different students in order to generalize in a best possible way and should therefore be user independent. The drawn gesture was mapped to a sensor array/lattice as shown in Figure 7 right. Each sensor node will be considered to be a motion node if the gesture has passed within the sensing range of the sensor. Each character has been drawn in a single-stroke gestural approach inspired by Graffiti [Hawkins 2002], Figure 8.

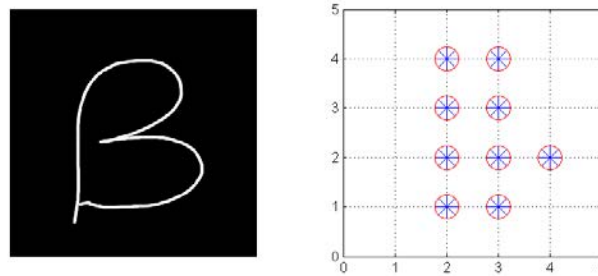


Fig. 7: Left: a template from the training set. right: resulting mapping of motion nodes in a 4x4 network



Fig. 8: Alphabet used in the Dataset

5.2. Algorithm Selection

Simulations assume an event-driven message exchange and emulate packet loss by randomly dropping packets with a certain likelihood. We have tested all four algorithms: Basic, Distributed Dataset (DD), Distributed Vector Computation (DV), and Hybrid (DDV) using data recorded from the iPad for a 10x10 network with feature vector length $N = 16$. We used the Bloom Filter-based multicasting algorithm from [Hosseinmardi et al. 2012] with $TTL = 2$ and threshold 10. Figure 9 shows the accuracy of all algorithms on a 10x10 network with vector length $N = 16$ without any packet loss. There is a small decrease (around 2%) due to the effect at the edges of the network where nodes do not have access to the complete data set via their neighbors.

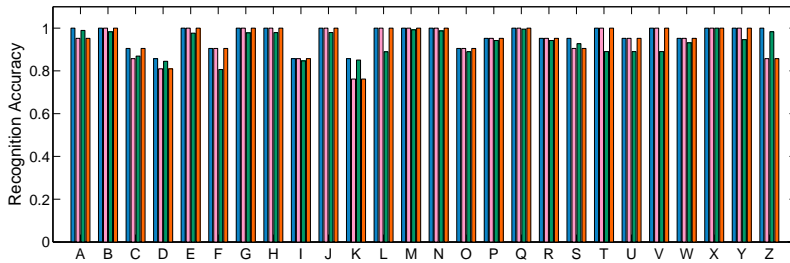


Fig. 9: Probability of gesture recognition for 546 templates of dataset for Basic (blue), DV (pink), DD (green), DDV (red) algorithm when there is no packet loss. Experiment on a 10 x 10 network, $N = 16$

Figure 9 shows that in the absence of packet loss, all the distributed algorithms work fairly well compared to the Basic algorithm, while also maintaining an advantage of memory and/or computation efficiency. They gain this property through sharing the data with their neighbors. Therefore we assessed the communication overhead of distributing packets in the advanced algorithms, as shown in Figure 10. As expected, distributing the computation as is done in both the DV and DDV algorithms results in greater communication overhead than merely distributing the dataset. However, the relative increase in overhead is mitigated by the fact that all algorithms including the Basic algorithm incur substantial communication overhead because of the initial sharing of sensor data via multicast packets with other group members.

We can use the algorithms with no special changes for a wireless environment prone to packet loss. Since all three distributed algorithms require data packet exchange to make a decision, we should consider the effect of packet loss on their recognition accuracy. Here, packet loss is simulated by dropping packets at a constant rate. In the following experiment we have considered the effect of packet loss. We test our algorithms on perfect input using an “A” shape without discontinuities. Figure 11 shows the detection probability as a function of packet loss probability and the chosen algorithm. We observe that the distributed vector computation approach is more resilient to packet loss than the distributed dataset approach, though even the distributed dataset algorithm accuracy only drops to 80% with 25% packet loss, which can be viewed as an acceptable recognition rate. The DDV algorithm stays fairly robust to packet loss, inheriting the robustness to loss of partial computation results of the DV algorithm and is an improvement over the DD algorithm. As a result of its superior memory and com-

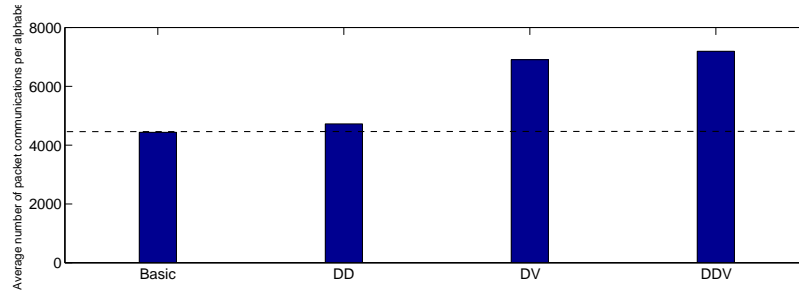


Fig. 10: Average number of packets exchanged for different algorithms (averaged over 546 inputs) when there is no packet loss. Experiment on a 10 x 10 network, $N = 16$

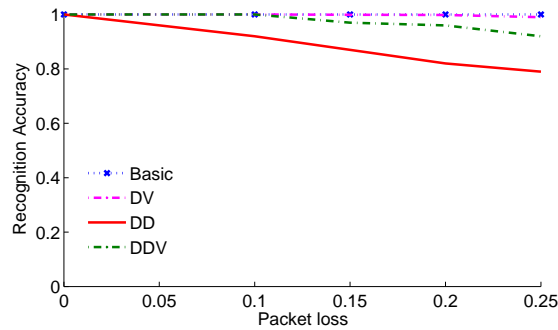
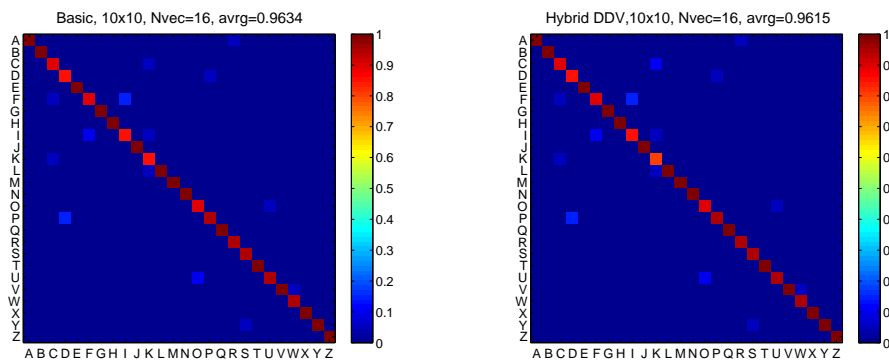


Fig. 11: Gesture recognition probability as a function of packet loss probability for an identical “A” shape compared to the Dataset shown in Figure 8. Experiment on a 10 x 10 network, $N = 16$

putational efficiency as well as resilience to packet loss, we choose to use the Hybrid DDV algorithm in our remaining analysis.



Another way to visualize the performance of the Hybrid DDV algorithm is via a confusion matrix, namely which gestures were identified versus which gestures were

Fig. 12: Confusion matrix of Basic (left) and Hybrid DDV Algorithm (right) when there is no packet loss. Experiment on a 10 x 10 network, $N = 16$

input, as shown in Figure 12. We observe that the Hybrid algorithm has a fairly similar confusion pattern compared to the Basic algorithm. DDV has the additional virtues of lower memory requirements and lower computation per node.

5.3. Sensitivity to Network Resolution and Vector Length

We next evaluate the effect of network size and feature vector length on algorithm performance. The dataset is comprised of 26 characters that have been written in a single stroke directional manner as has been shown in Figure 8. We have distributed these 26 characters between neighbors such that each node can have access to the whole alphabet via its neighbors. As we mentioned in section 3.1 a gesture will be considered as a chain vector with fixed length N . To examine the role of this parameter in the recognition algorithm, we have evaluated the recognition accuracy for $N = 6, 8, 10, 16, 25, 32, 64$ across a range of networks (4x4, 6x6, 10x10 and 12x12), as shown in Figure 13.

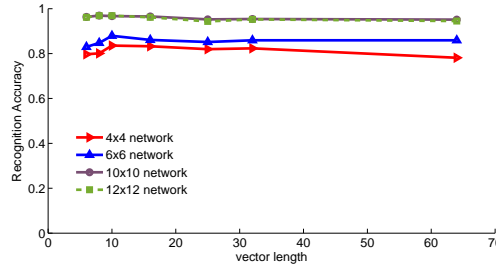


Fig. 13: Average recognition accuracy of Hybrid Algorithm when there is no packet loss. Experiment on four different network sizes, $N = 6, 8, 10, 16, 25, 32, 64$

As we expected the recognition accuracy was reduced as the network size decreased. This happened because the drawn gesture was captured by fewer sensors and as a result the input character has lower resolution. For each network size, as we change the feature vector length we see that there is a very slight change of maximum 4.9% in the recognition accuracy. On the other hand, as the network size increases, the number of exchanged packets increases, as seen in Figure 14.

In Figure 15, we look more closely at the effect of changing vector length on each alphabet recognition rate in a 10x10 network. We can see that increasing or decreasing vector length does not have the same effect on recognition for different characters. For example, 'D' has the best result at $N = 10$, while 'U' achieves the best recognition at $N = 16$. Surprisingly, for all characters, the recognition rate using length $N = 32$ does not exceed the rate when either $N = 16$ or $N = 10$.

In Figure 16, for a network size of 4x4 we also typically observe the worst result for $N = 32$. However, also note that there are some characters like 'Q', 'S' and 'W' for which the accuracy has the highest rate at $N = 32$.

Figure 17 shows that reducing network size, which causes reduction on input character resolution, does not always have a negative effect on recognition (however this is mostly true). For example looking closely we observe that for 'C' and 'I' recognition

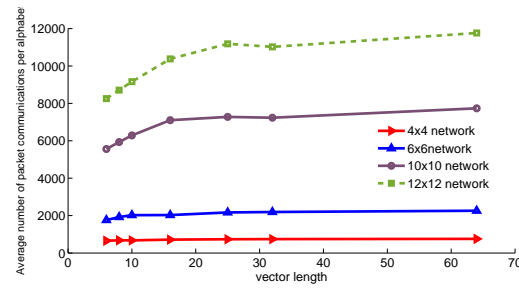


Fig. 14: Average number of packets exchanged (averaged over 546 input) using Hybrid Algorithm when there is no packet loss. Experiment on four different network sizes, and $N = 6, 8, 10, 16, 25, 32, 64$.

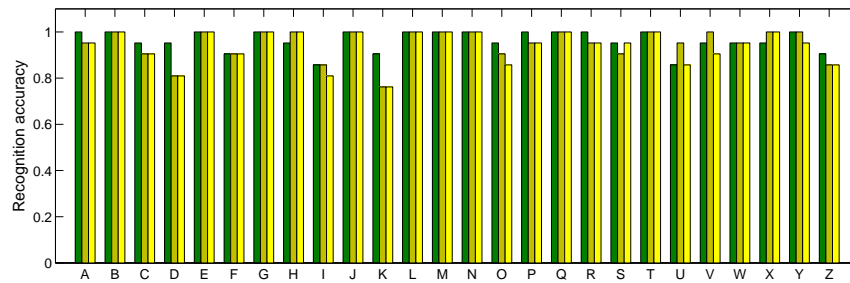


Fig. 15: Average recognition accuracy of Hybrid Algorithm for 10x10 network and chain vector length N equal to 10 (dark green), 16 (light green), 32 (yellow)

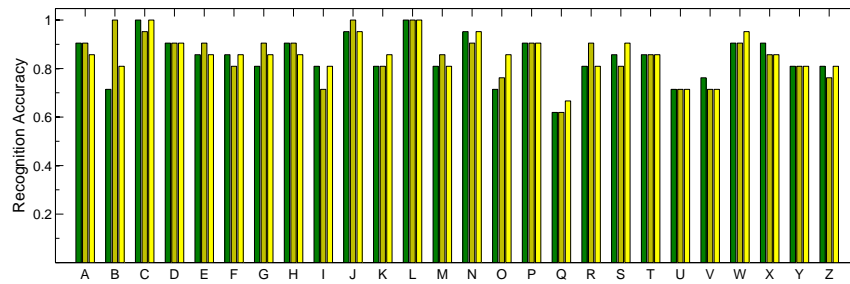


Fig. 16: Average recognition accuracy of Hybrid Algorithm for 4x4 network and chain vector length N equal to 10 (dark green), 16 (light green), 32 (yellow)

accuracy has increased. On the other hand, we can see that characters like 'M', 'Q' and 'Y' are worst affected by reducing the network size.

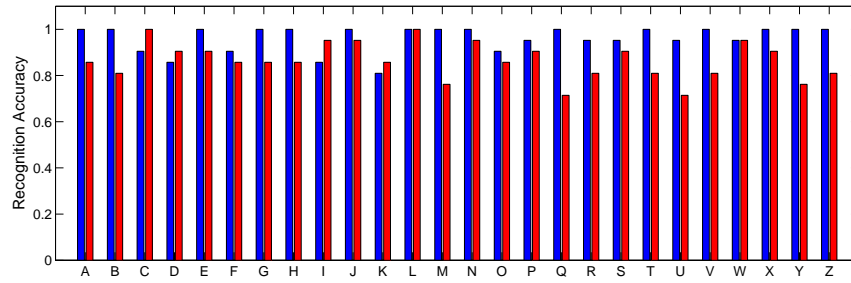


Fig. 17: Average recognition accuracy of Hybrid Algorithm for 10x10 (blue) and 4x4 (red) network, $N = 16$

5.4. Selecting a Subset of an Alphabet

In some applications a sensor lattice may be required to respond to only a small number of gestures. One research goal is therefore to choose a subset of gestures that are most widely differentiated so that it is easiest to tell them apart. In this section we show how we can choose a specific number of characters from the whole 26-character alphabet to achieve the best recognition percentage. As a case study, suppose we need to choose 5 characters, where the number five is chosen to allow us to optimally distribute the dataset on a 4-neighborhood grid structure so that each node has access to the entire dataset via its one hop neighbors. In the following, we perform a heuristic experiment considering all $\binom{26}{5} = 65780$ possible combinations of 5-member groups. Figure 18 shows the number of 5-member groups that achieve an accuracy of 100% for different network sizes and vector length. We can see from Figure 18 that the worst case is for a network size 4x4 and vector length $N = 32$, though there still exist more than 2500 different combinations which achieve full recognition.

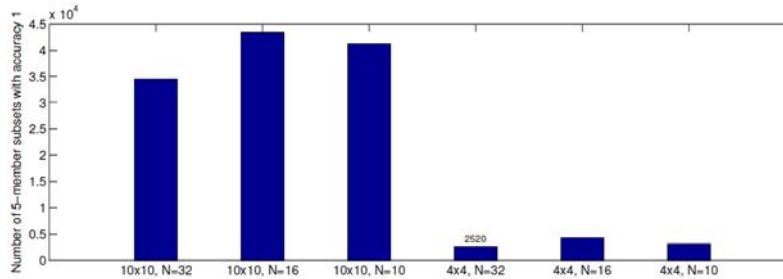


Fig. 18: Number of 5-member groups which have average accuracy 1

The question that arises here is how we can choose the subset that maximizes accuracy, containing the most mutually distant characters in the feature space. We can introduce an optimization problem to find the characters which have the most distinct feature vectors. The following experiment shows how different combinations act differently in the presence of packet loss. Figure 19 shows the difference between choosing two 5-member sets. In one group, we attempt to distinguish B from the set of “A, B, D, H, N” of which four of the members have an upward start. In the other group “B, H, O, W, Z”, each of the members have different starting directions. Both groups have

the recognition accuracy 100% in the absence of packet loss. We observe how the group “A, B, D, H, N” is more sensitive to packet loss than group “B, H, O, W, Z”. This example illustrates that other factors such as differences in the initial directionality of the characters may be helpful in selecting a subset of the alphabet that maximizes accuracy. More exploration is required to determine additional factors like starting direction that would be useful in influencing optimal subset selection.

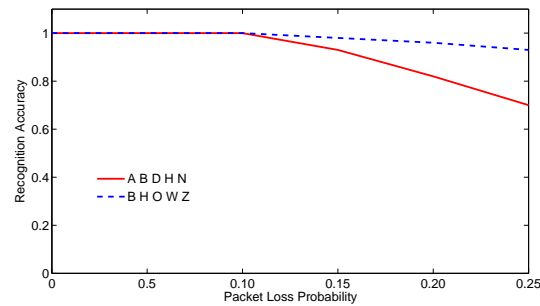


Fig. 19: Probability of gesture recognition for an identical “B” shape compared to two different 5 members set, Network size 4x4, N=16, Hybrid DDV algorithm

6. DISCUSSION

All experiments, including those for data gathered from our hardware setup, were simulated on a desktop computer. This allowed us to explore a wide range of memory and bandwidth requirements, not necessarily supported by a particular hardware setup. Whereas all of the presented algorithms can be run on the hardware used here for the type of data presented (the hybrid algorithm exchanges on average 2490 packets, each a few bytes in length), this might not be the case for larger datasets or gestures with higher resolutions, which will significantly slow down recognition and/or reduce accuracy if bandwidth is low.

Simulation results confirm that the distributed dataset approach is not very robust to packet loss despite being memory efficient. Distributed computation, in which only parts of the gesture are lost in the event of packet loss, turns out to be significantly more robust, but scales poorly with increasing size of the dataset. The hybrid approach maintains the memory efficiency of the distributed datasets, while combining the computational efficiency of both. Performance in terms of recognition probability, however, is identical to that of the distributed vector computation approach until packet loss is above around 10% for a 10x10 network with 26 different characters in the dataset.

The primary cost of the hybrid DDV algorithm is increased packet overhead compared to the Basic algorithm, though that is mitigated somewhat by the initial costs required in the Basic algorithm to share packets of raw sensor data among the motion nodes. We believe that many sensor array systems may find the additional cost of packet communication reasonable compared to the savings in memory and computation achieved by the hybrid DDV algorithm, especially those systems where power to the sensor lattice is not an issue, such as in smart sensor walls.

A common challenge in recognition of written characters is that sub-sampling bears the risk of picking a vector that is not oriented toward the overall motion of the gesture, e.g. due to jitter of the user’s hand. While this problem could be prevented by

an additional smoothing step before sub-sampling, we believe that jitter is a problem more specific to hand writing than to larger-scale hand motions. Yet, our results show that low resolution negatively affects pattern recognition.

Distributing the dataset exclusively among neighboring nodes allows minimizing the communication overhead of a distributed approach. By this, the overall size of the dataset is limited by the available memory on every node and the communication topology. In this paper, we consider four and six member neighborhoods, but also 8-member neighborhoods are possible, depending on the available hardware. We also note that higher connectivity degrees might increase the detection rate in presence of packet loss for the hybrid approach.

We note that our current approach is not invariant to rotations. Doing this would require an additional normalization step that transforms each chain vector into a rotation-invariant representation. We note, however, that this would also limit the system to rotation-invariant gestures, possibly excluding e.g., “C” and “U”.

7. CONCLUSIONS

We presented a suite of distributed gesture recognition algorithms for next generation smart materials, which embed minimalist sensing, computation and communication, and allow a trade off of computation and memory for more communication among nodes. All algorithms are fully distributed, scalable to the size of the sensing lattice, and could be used on amorphous structures, that is structures covered with the sensor lattice which have arbitrary outer shape and possibly holes in the surface. Algorithms have been validated and tested with respect to packet loss in simulation using digitized hand-drawn gestures. These results show that a hybrid approach, which distributes both computation and the template dataset among nodes, can lead to high detection probabilities when packet loss is less than 10% on a grid lattice topology, even if the input gesture is presented with discontinuities due to either sensor failure or packet loss.

In the future, we wish to extend our work to multi-stroke alphabets as well as multi-alphabet input commands. We are also interested in validating the proposed advanced algorithm on physical hardware, and are in the process of building a smart wall system. We plan to implement Bloom-filter based multi-cast communication and a basic character recognition algorithm on the hardware shown in Figure 1, and are also interested to implement the proposed algorithms on other smart materials such as skins and shape changing materials.

REFERENCES

- C. Bahlmann, B. Haasdonk, and H. Burkhardt. 2002. Online handwriting recognition with support vector machines—a kernel approach. In *Frontiers in Handwriting Recognition, 2002. Proceedings. Eighth International Workshop on*. IEEE, 49–54.
- W. Butera. 2002. *Programing a paintable computer*. Ph.D. Dissertation. Ph. D. dissertation, Program in Media Arts and Sciences, School of Architecture and Planning, MIT.
- J. M. Carmona and J. Climent. 2012. A Performance Evaluation of HMM and DTW for Gesture Recognition (*Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*), Luis Alvarez, Marta Mejail, Lus Gmez Dniz, and Julio C. Jacobo (Eds.). Springer, 236–243.
- W.J. Confer and R.O. Chapman. 2004. System and method of handwritten character recognition. (April 13 2004). US Patent 6,721,452.
- RH Davis and J. Lyall. 1986. Recognition of handwritten characters—a review. *Image and Vision Computing* 4, 4 (1986), 208–218.
- M. Elmezain, A. Al-hamadi, and B. Michaelis. 2009. A Hidden Markov Model-Based Isolated and Meaningful Hand Gesture Recognition.

- N. Farrow, N. Sivagnanadasan, and N. Correll. 2014. Gesture Based Distributed User Interaction System for a Reconfigurable Self-Organizing Smart Wall. In *Proceedings of the 8th International Conference on Tangible, Embedded and Embodied Interaction (TEI)*. ACM, 245–246.
- J. Hawkins. 2002. (2002). [http://en.wikipedia.org/wiki/Graffiti\(Palm.OS\)](http://en.wikipedia.org/wiki/Graffiti(Palm.OS))
- H. Hosseinmardi, N. Correll, and R. Han. 2012. Bloom Filter-Based Ad Hoc Multicast Communication in Cyber-Physical Systems and Computational Materials. *Wireless Algorithms, Systems, and Applications* (2012), 595–606.
- D. Hughes and N. Correll. 2014. A Soft, Amorphous Skin that can Sense and Localize Texture. In *IEEE International Conference on Robotics and Automation (ICRA)* (2014-05-31). Hong Kong, 1844–1851. <http://ieeexplore.ieee.org/xpl/articleDetails.jsp?tp=&arnumber=6907101>
- H.K. Lee and J.H. Kim. 1999. An HMM-based threshold model approach for gesture recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 21, 10 (1999), 961–973.
- J. Lifton, D. Seetharam, M. Broxton, and J. Paradiso. 2002. Pushpin computing system overview: A platform for distributed, embedded, ubiquitous sensor networks. *Pervasive Computing* (2002), 139–151.
- V. J. Lumelsky, M. S. Shur, and S. Wagner. 2001. Sensitive Skin. *IEEE Sensors Journal* 1, 1 (2001), 41–51.
- S. Ma, H. Hosseinmardi, N. Farrow, R. Han, and N. Correll. 2012. Establishing Multi-Cast Groups in Computational Robotic Materials. In *IEEE International Conference on Cyber, Physical and Social Computing* (2012-11-20). Besancon, France.
- M. A. McEvoy and N. Correll. 2014. Thermoplastic variable stiffness composites with embedded, networked sensing, actuation, and control. *Journal of Composite Materials* (2014). <http://jcm.sagepub.com/content/early/2014/03/07/0021998314525982>
- M. A. McEvoy and N. Correll. 2015. Materials that couple sensing, actuation, computation and communication. *Science* 347, 6228 (2015). <http://www.sciencemag.org/lookup/doi/10.1126/science.1261689>
- M. A. McEvoy, N. Farrow, and N. Correll. 2013. Robotic Materials with Controllable Stiffness. In *Proceedings of the 19th International Conference on Composite Materials (ICCM)*.
- S. Mitra and T. Acharya. 2007. Gesture recognition: A survey. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on* 37, 3 (2007), 311–324.
- K. Murakami and H. Taguchi. 1991. Gesture recognition using recurrent neural networks. In *Proceedings of the SIGCHI conference on Human factors in computing systems: Reaching through technology*. ACM, 237–242.
- J. Nagi, H. Ngo, A. Giusti, Luca M. Gambardella, J. Schmidhuber, and G. A. Di Caro. 2012. Incremental learning using partial feedback for gesture-based human-swarm interaction.. In *RO-MAN*. IEEE, 898–905.
- Z. Nakad, M. Jones, and T. Martin. 2003. Communications in Electronic Textile Systems. In *Proceedings of the 2003 International Conference on Communications in Computing*. 37–43.
- O.F. Ozer, O. Ozun, C.O. Tuzel, V. Atalay, and A.E. Cetin. 2001. Vision-based single-stroke character recognition for wearable computing. *Intelligent Systems, IEEE* 16, 3 (2001), 33–37.
- J. A Paradiso, J. Lifton, and M. Broxton. 2004. Sensate mediamultimodal electronic skins as dense sensor networks. *BT Technology Journal* 22, 4 (2004), 32–44.
- Halley Profita, Nicholas Farrow, and Nikolaus Correll. 2015. Flutter: An Exploration of an Assistive Garment Using Distributed Sensing, Computation and Actuation. In *Proceedings of the Ninth International Conference on Tangible, Embedded, and Embodied Interaction (TEI '15)*. ACM, New York, NY, USA, 359–362. DOI : <http://dx.doi.org/10.1145/2677199.2680586>
- H. Sakoe and S. Chiba. 1990. Readings in Speech Recognition. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, Chapter Dynamic Programming Algorithm Optimization for Spoken Word Recognition, 159–165. <http://dl.acm.org/citation.cfm?id=108235.108244>
- T. Schlömer, B. Poppinga, N. Henze, and S. Boll. 2008. Gesture recognition with a Wii controller. In *Proceedings of the 2nd International Conference on Tangible and Embedded Interaction*. ACM, 11–14.
- A.D. Wilson and A.F. Bobick. 1999. Parametric hidden markov models for gesture recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 21, 9 (1999), 884–900.
- Y. Wu and T. Huang. 1999. Vision-based gesture recognition: A review. *Gesture-based communication in human-computer interaction* (1999), 103–115.