

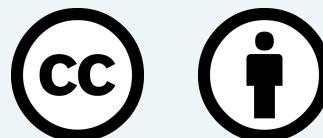
On-line Gesture Recognition

Luis A. Leiva

luileito@prhlt.upv.es

PRHLT Research Center

Departamento de Sistemas Informáticos y Computación
Universitat Politècnica de València



<http://creativecommons.org/licenses/by/4.0/>

Presentation Outline

Introduction	1
Preliminaries	19
Some Techniques	31
Recap	59
References	63

Slides available at <http://personales.upv.es/luileito/lectures/gestures-upv.pdf>

Introduction

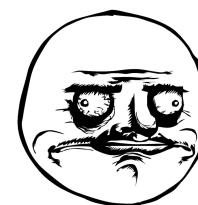
Definition

Gesture /'dʒɛstʃə/ noun

The use of motions of the limbs or body as a means of expression.

synonyms: signal, sign, motion, indication, gesticulation

“Gestures are hand-drawn strokes that do things.”



— Lipscomb (1991)

Definition

- **Off-line** gesture recognition:
 - post-hoc*, processed after user interaction
 - static data, no temporal info available
- **On-line** gesture recognition:
 - realtime, direct manipulation
 - sequential, time series data



Historical Precedents

sketchpad



Sutherland (1963)

RAND tablet



Davis and Ellis (1964)

Gestures Today



Minority Report. Image by 20th Century Fox & DreamWorks

Input Devices

Wii. Image by Nintendo Co., Ltd.



Input Devices



T(ether). Image by Massachusetts Institute of Technology

Input Devices

Kinect for Xbox 360. Image by Microsoft Corporation



Input Devices



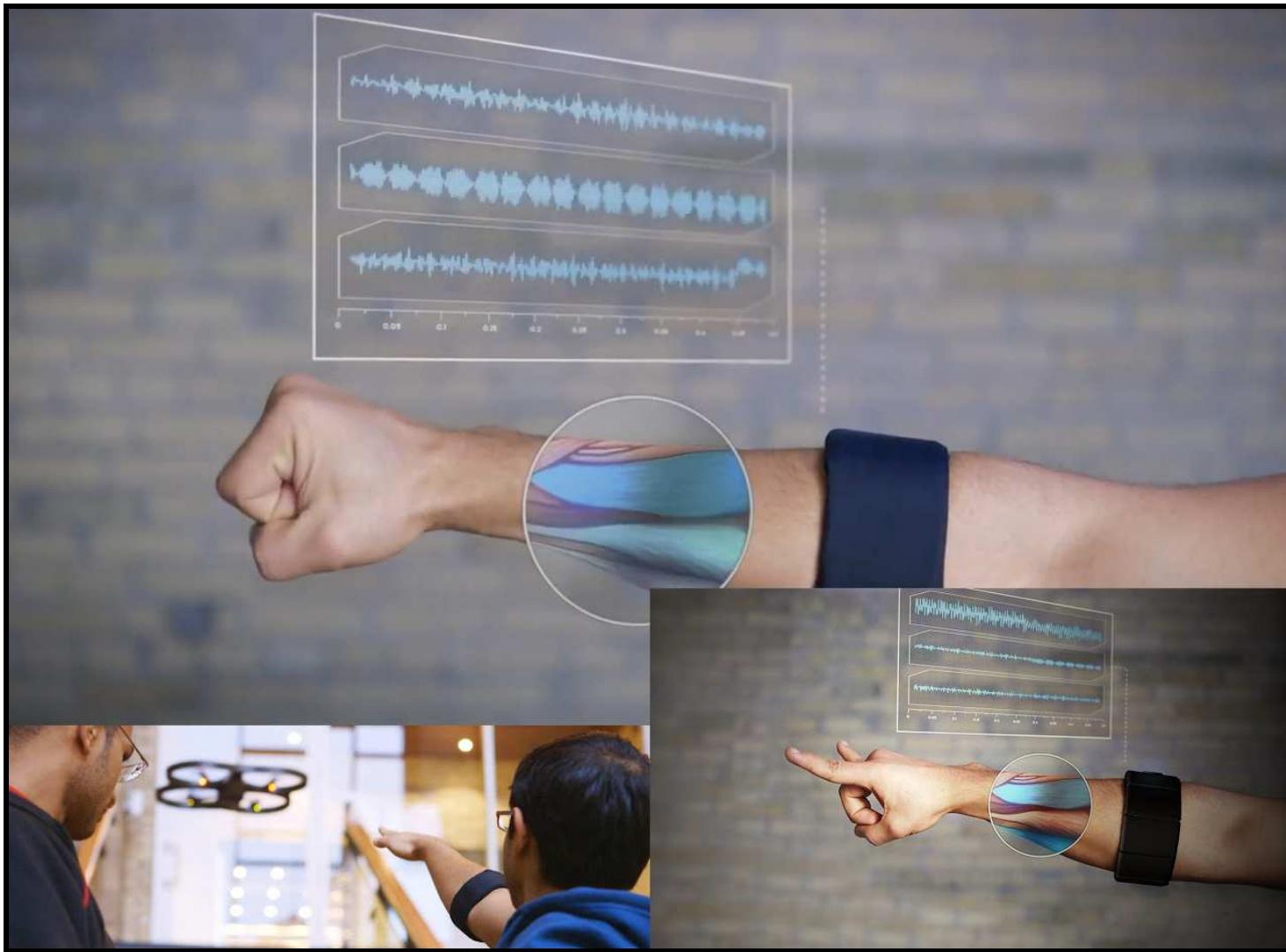
Humantenna. Image by Microsoft Research

Input Devices



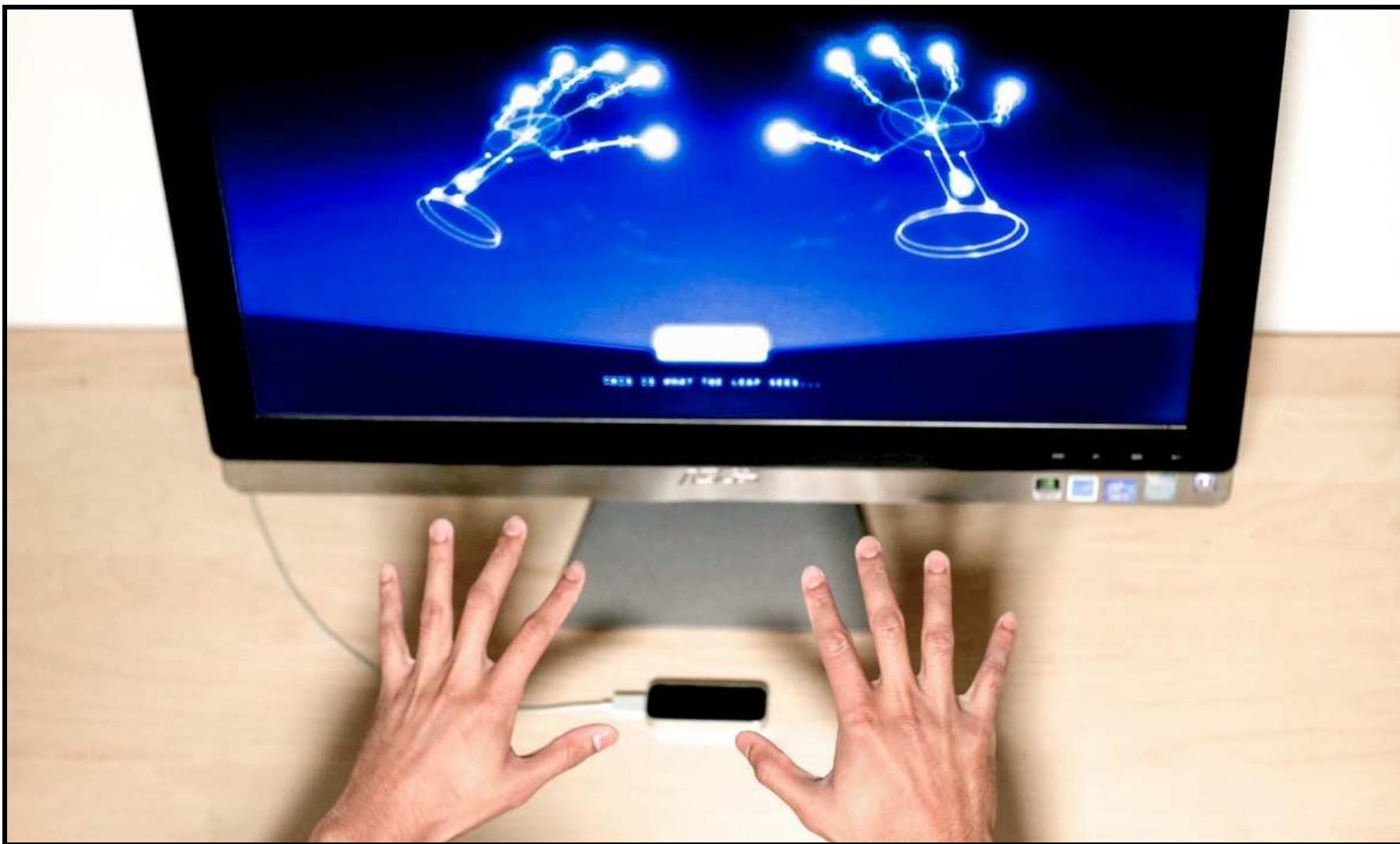
Skinput. Image by Carnegie Mellon University

Input Devices



Myo. Image by Thalmic Labs

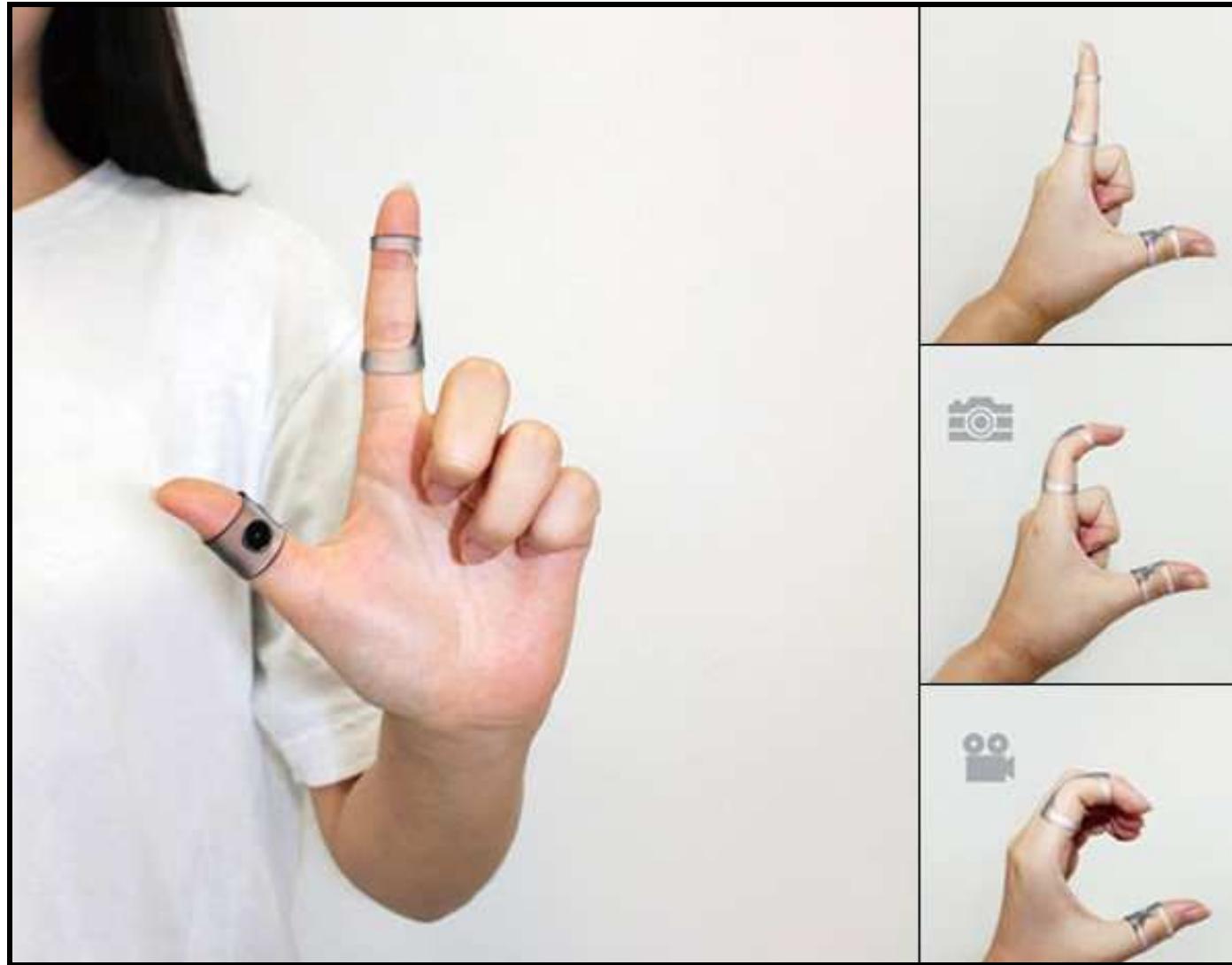
Input Devices



Leap motion. Image by Leap Motion, Inc.

Input Devices

Air Clicker. Image by Yanko Design

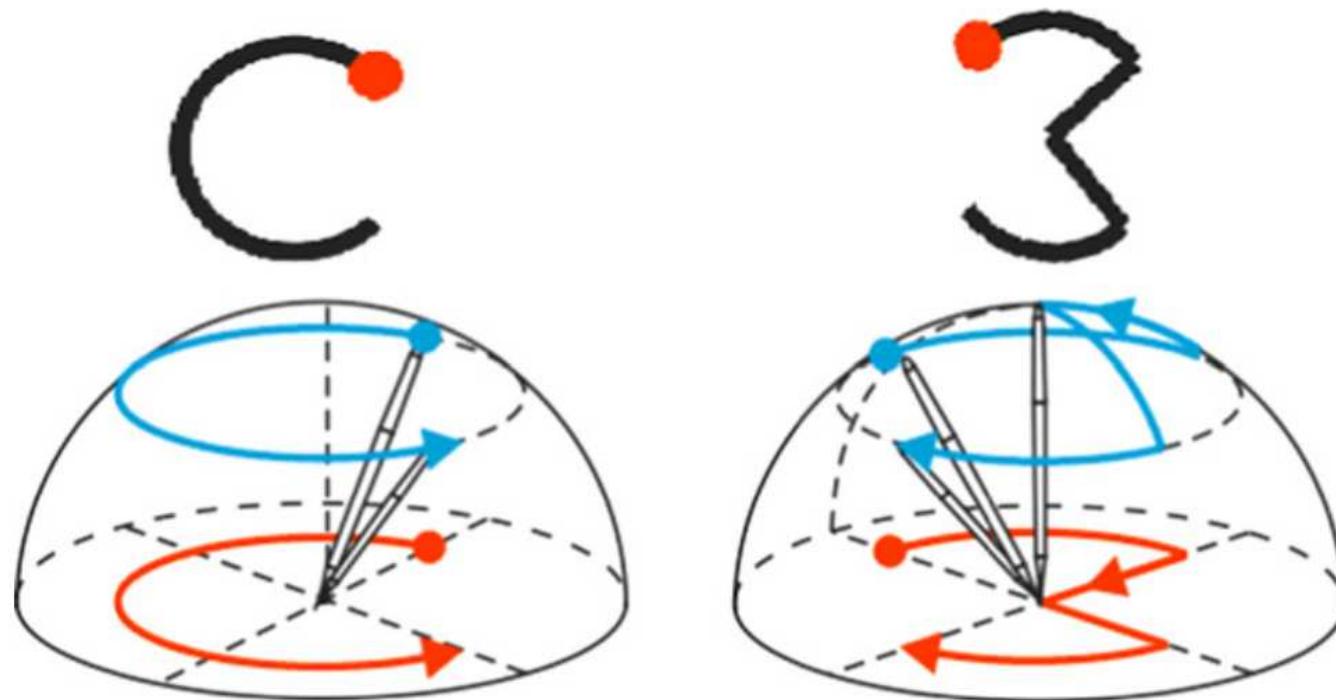


Input Devices



TapTap. Image by Woodenshark LLC

Input Devices



Pen Tail gestures, by Tian et al. (2012)

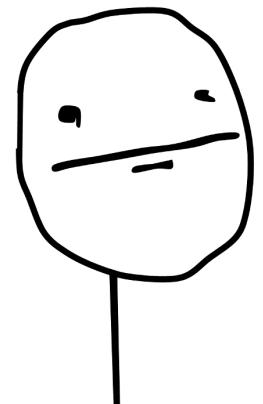
Advantages

- Natural communication
- Expressiveness
- Ergonomics
- Usability
- Fun

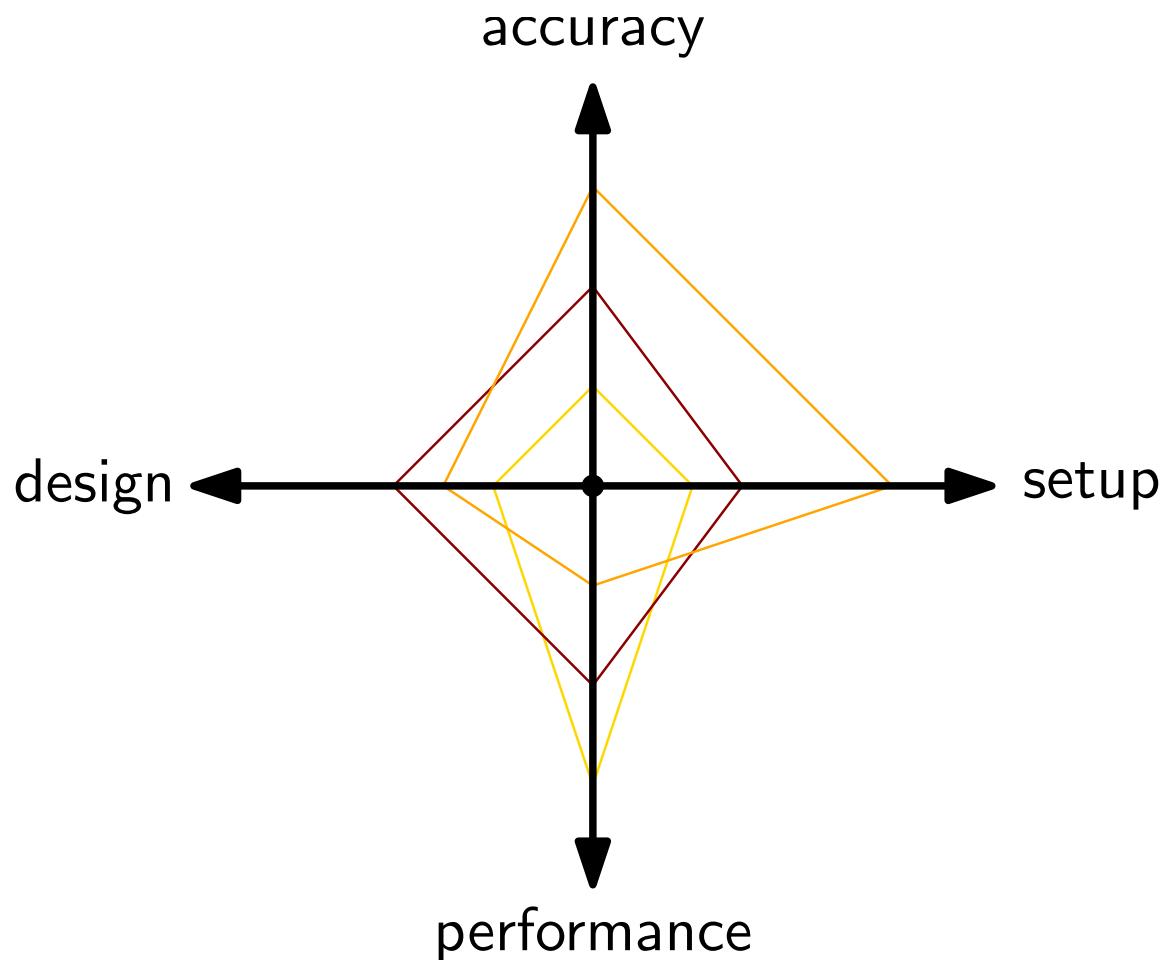


Disadvantages

- May break fundamental interaction principles:
Discoverability, **Reliability**, Scalability, etc.
- Ambiguity: **non-deterministic decoding**
- Lack of standards
- Cultural issues



Trade-offs



Preliminaries

Interaction Paradigms

Mid-air

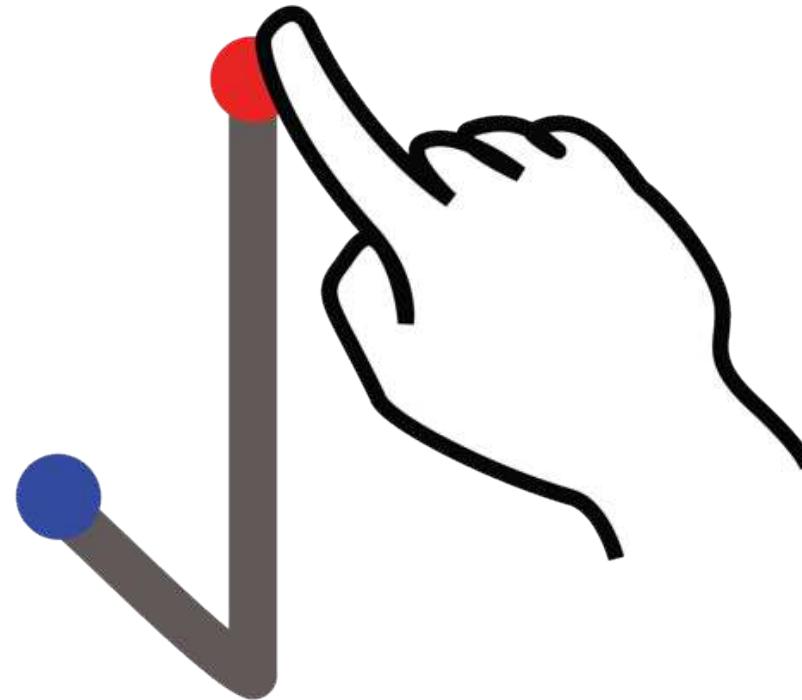


Onscreen



Definition

stroke = pointer **down** → pointer **move** → pointer **up**



$$s = \{(x_1, y_1, t_1) \dots (x_j, y_j, t_j) \dots (x_N, y_N, t_N)\}$$

A Taxonomy

zero-order gestures

one finger tap



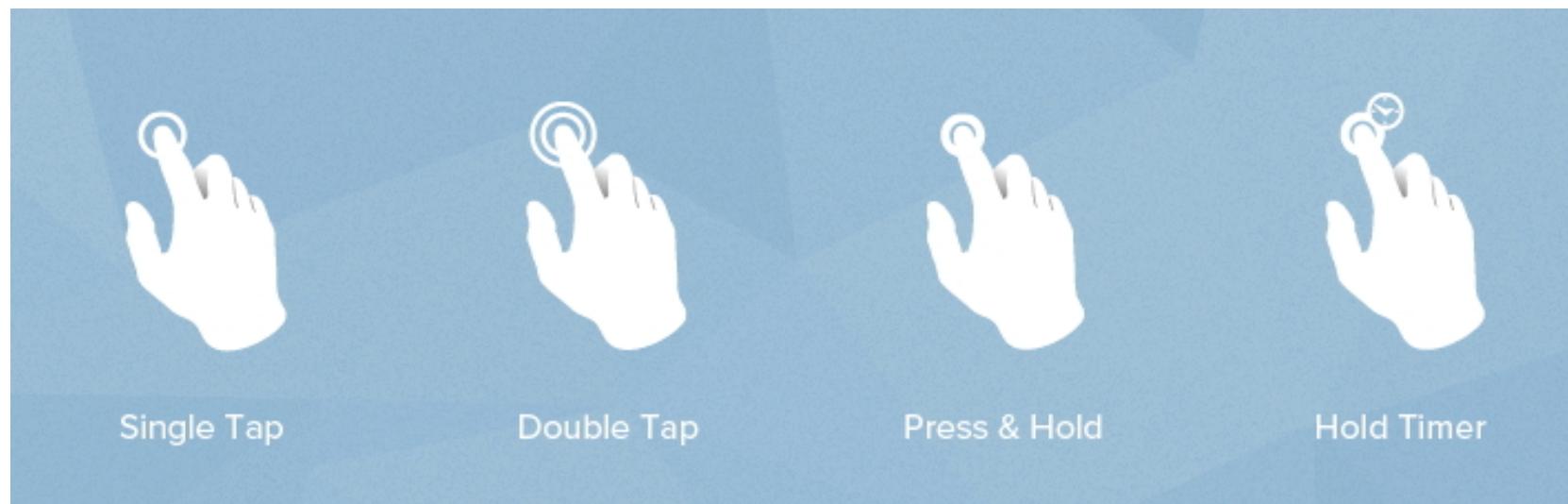
two finger tap



one finger double tap

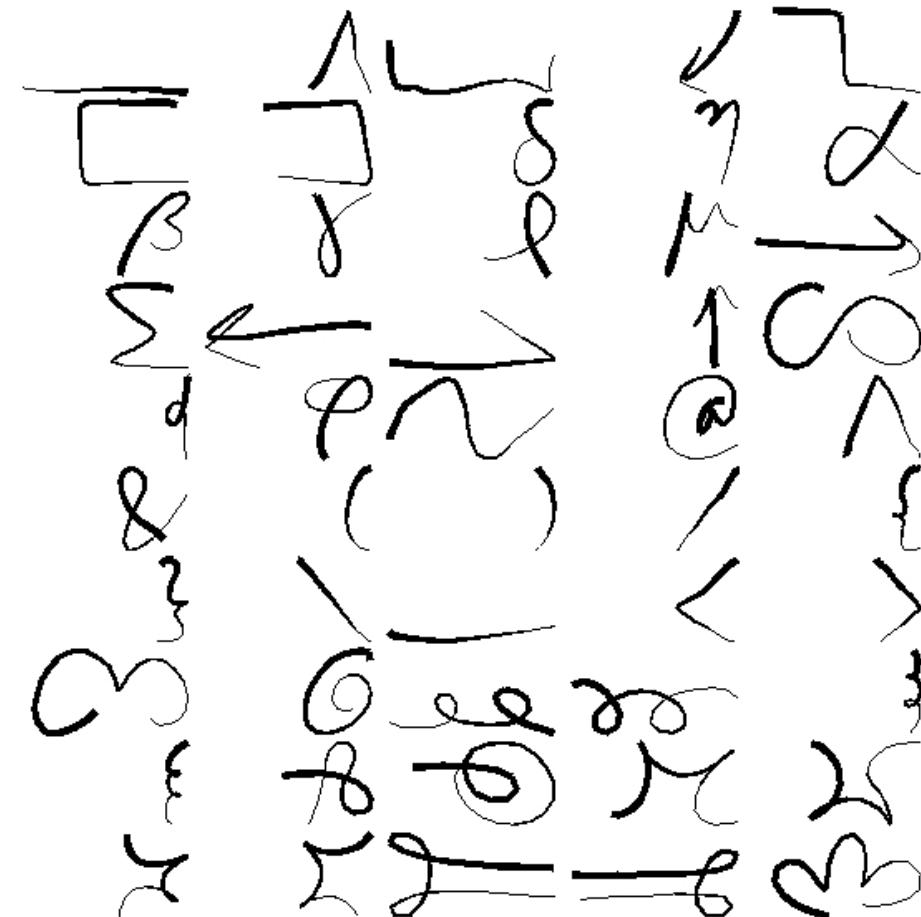
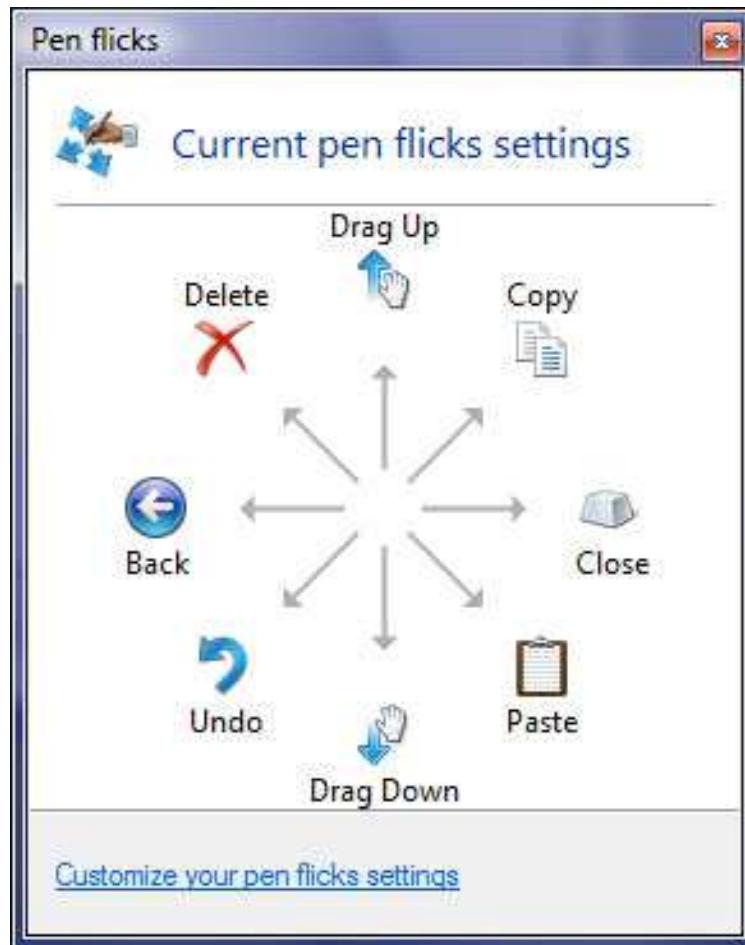


two finger double tap



A Taxonomy

first-order gestures (unistrokes)

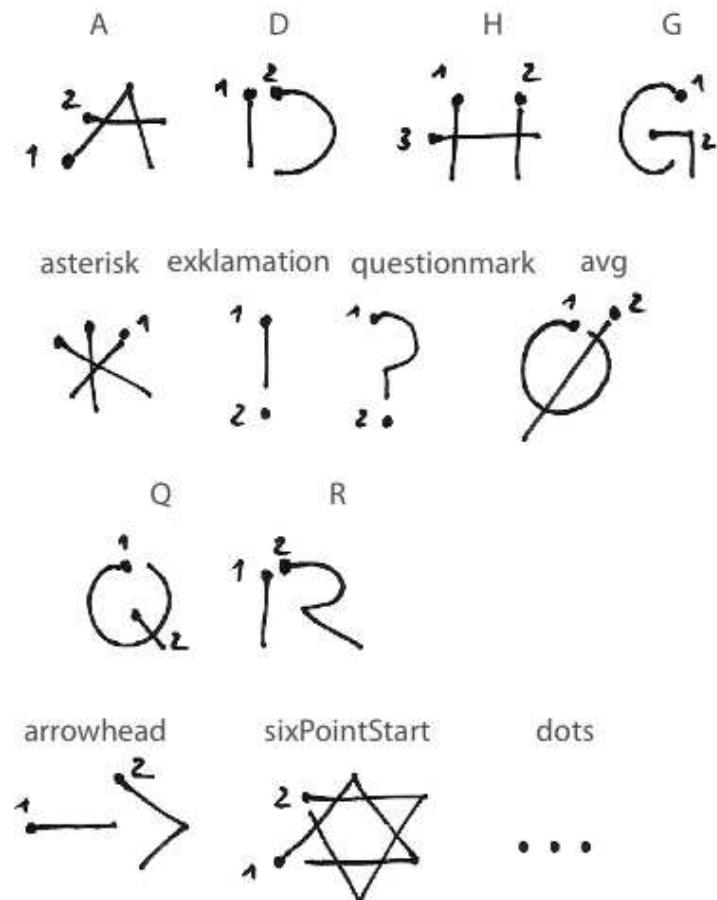


A Taxonomy

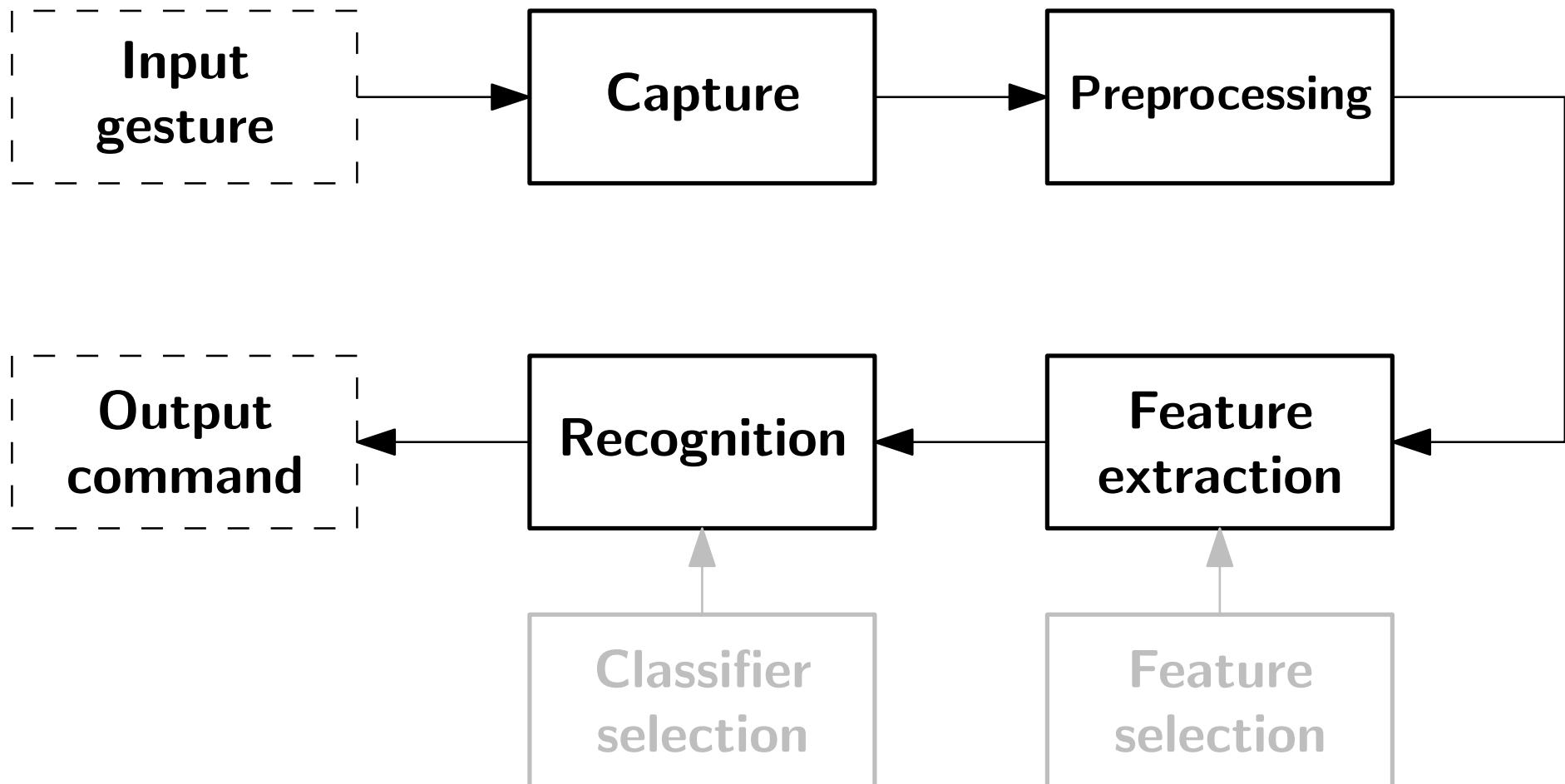
higher-order gestures (multistrokes)

IN MARGIN	IN TEXT
	insert word or letter
	delete, delete and close up space
	close up space
#	insert space
	equalize space; make space between words or lines equal
	begin new paragraph or continue last paragraph
	center
	flush left
	flush right
	reverses the order; transpose
	ragged margin; don't justify lines
	move text down
	move text up
	superscript 1 or subscript 2 (πr^2 or H_2O)
	spell out (set 1 hr as one hour)
	don't change; go back to the original
	change from Capital to lowercase letter (capital)

IN MARGIN	IN TEXT
	set in small capital letters (SMALL CAPITAL LETTERS)
	change from lowercase to capital (Capital)
	set in italic or slanted type (italic)
	set in Roman type (Roman)
	set in boldface type (boldface)
	wrong front or type style or size; set in correct type (correct type)
	insert comma
	insert period or colon
	insert double quotation marks (The Catbird Seat)
	insert single quotation mark or apostrophe (today's newspaper)
	insert hyphen (first class)
	insert en dash (3-4 credits)
	insert em dash (required courses- stand-alones or clusters)
	insert question mark (Who's on first)
	insert equals sign (1+1=2)
	insert parentheses or square brackets

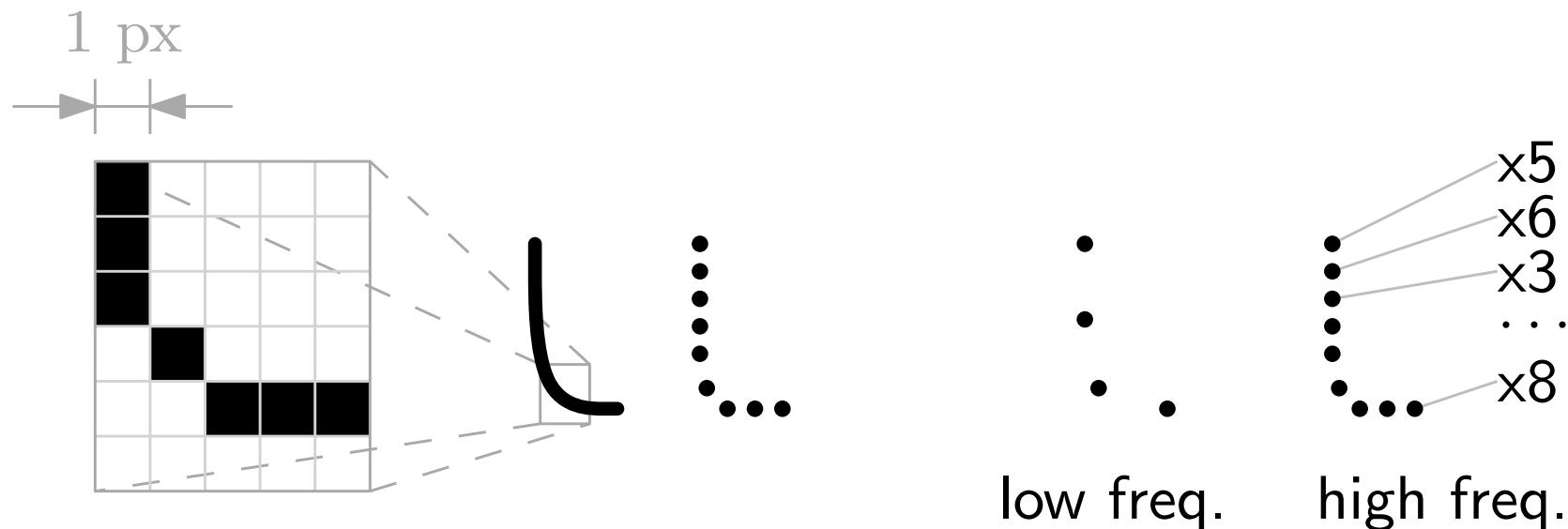


Processing Pipeline



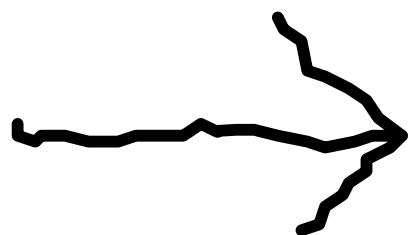
Capture

- Event-based
- Polling (constant frequency)

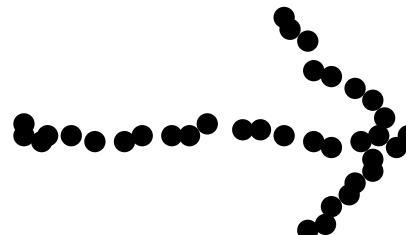


Sampling rate matters!

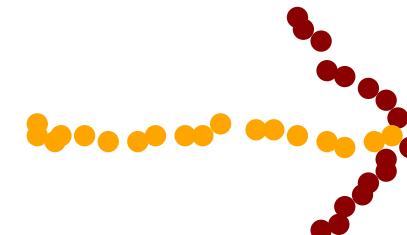
Preprocessing



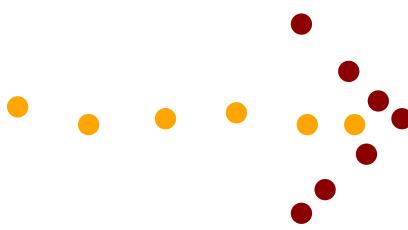
input



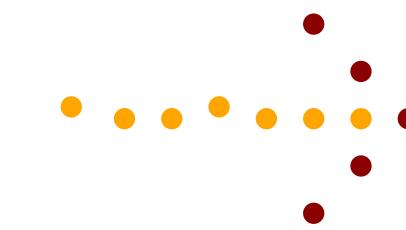
capture



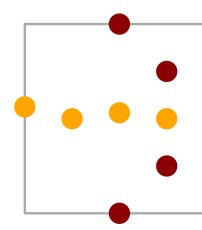
segmentation



noise removal*



resampling*



normalization*

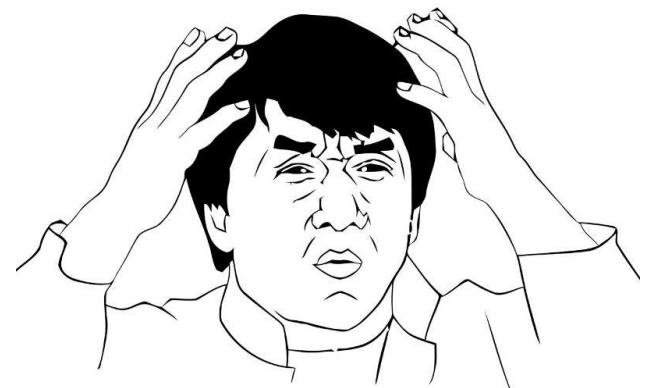
*optional steps

Features

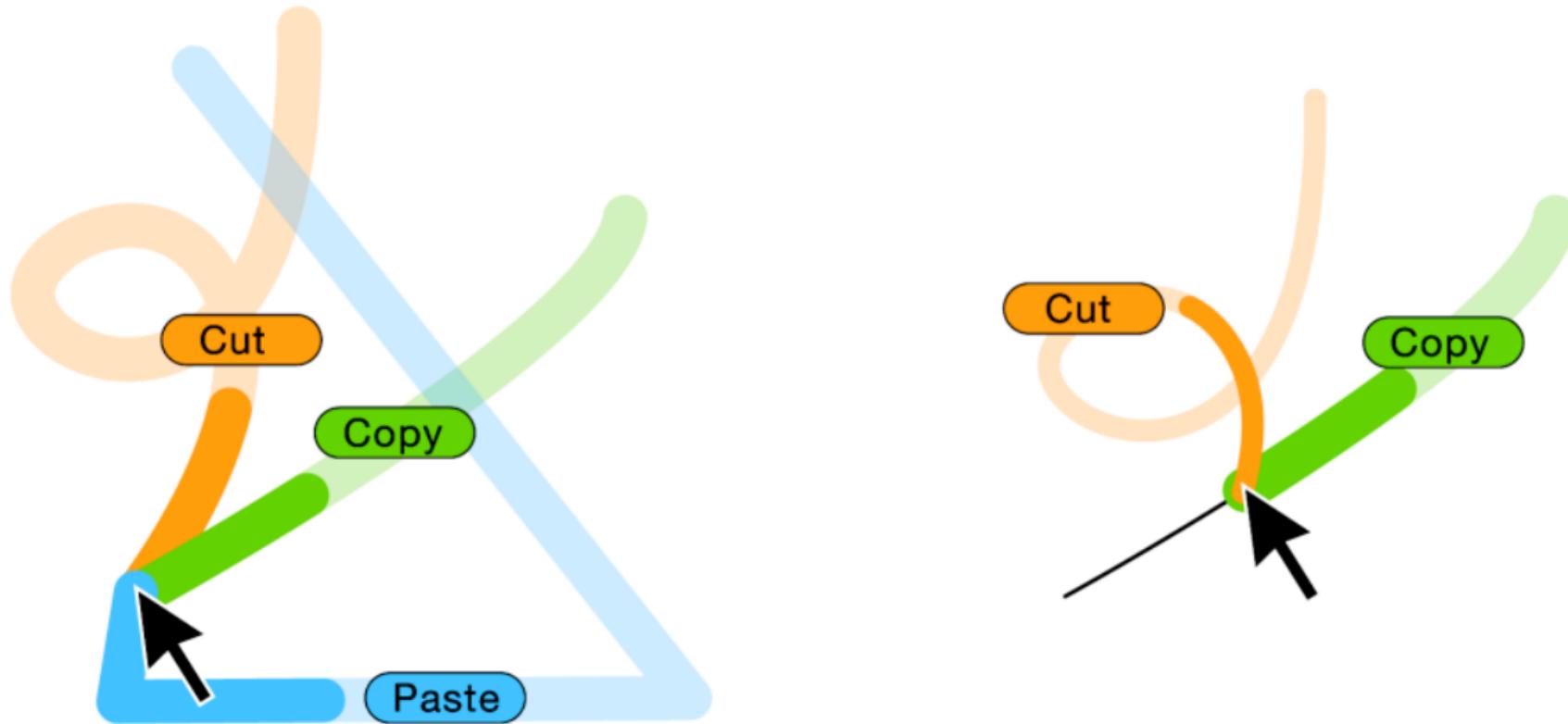


Recognition Techniques

- Hashing: dictionary lookup, zone coding, chain codes
- Parametric: linear fitting, corner detection
- Matching: DTW, k-NN, “dollar family”
- Statistical: linear classifier, NNs, HMMs, CRFs
- Ad hoc: knowledge-based, decision trees, FSM



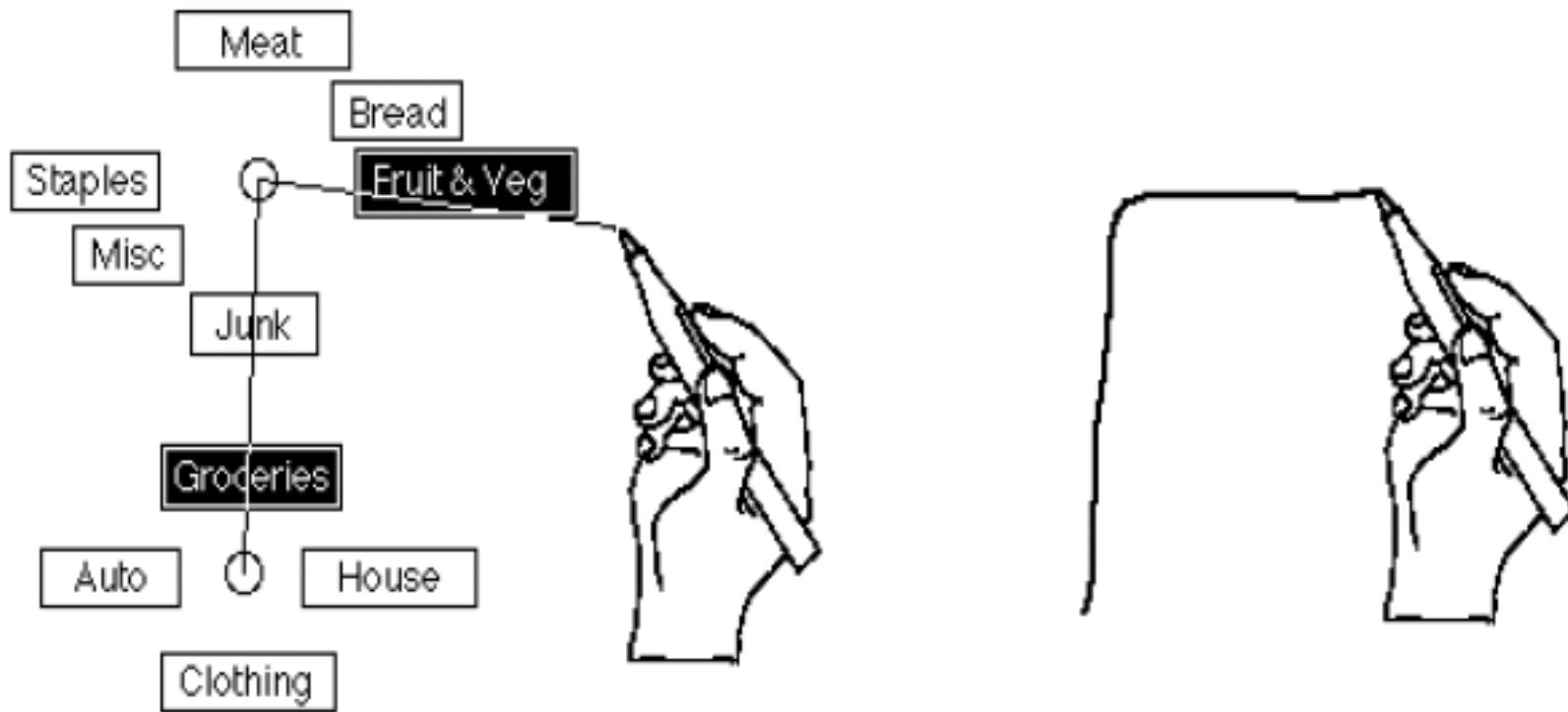
Continuous Recognition



OctoPocus, by Bau and Mackay (2008)

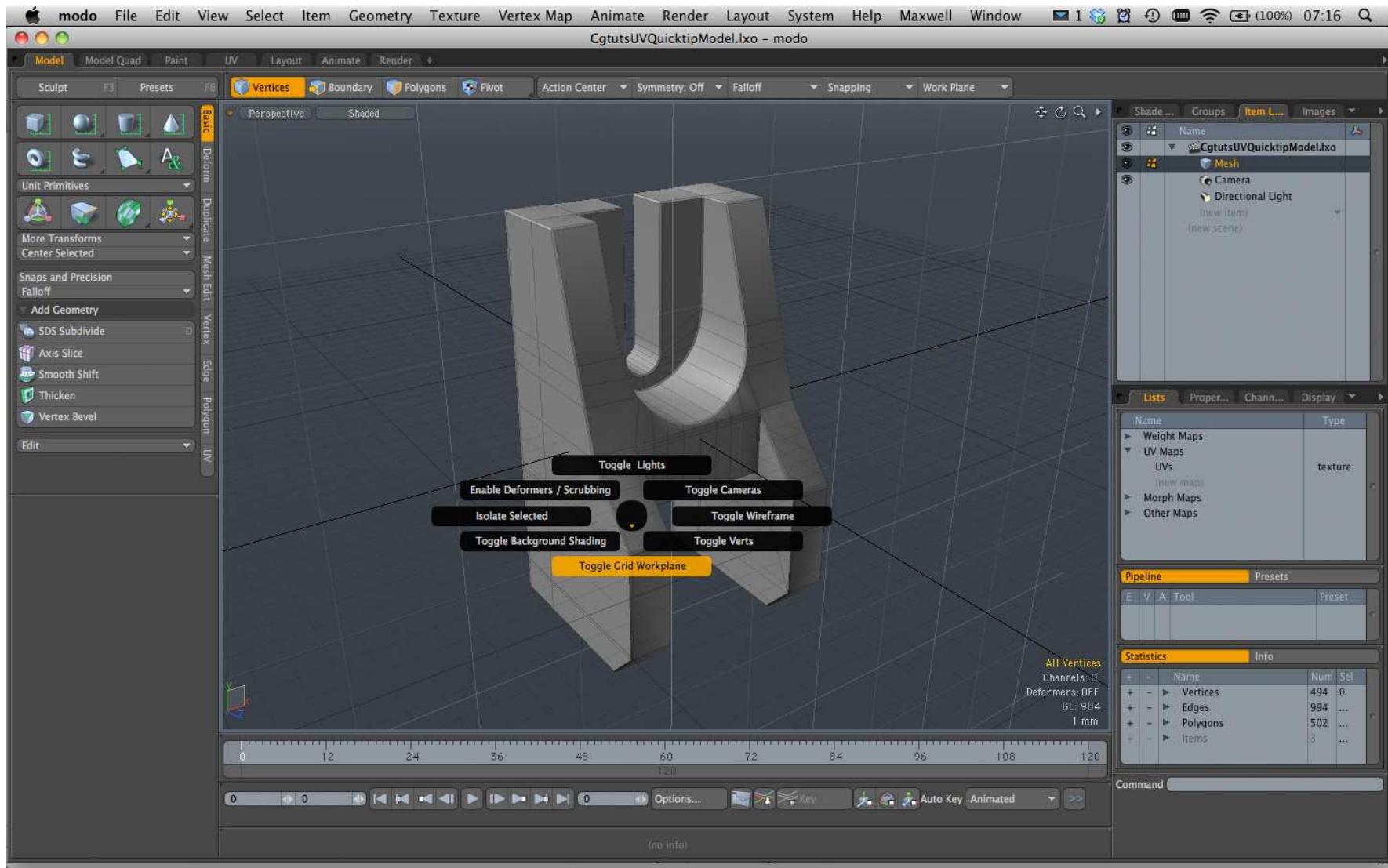
Some Techniques

Marking Menus



by Kurtenbach (1991)

Marking Menus



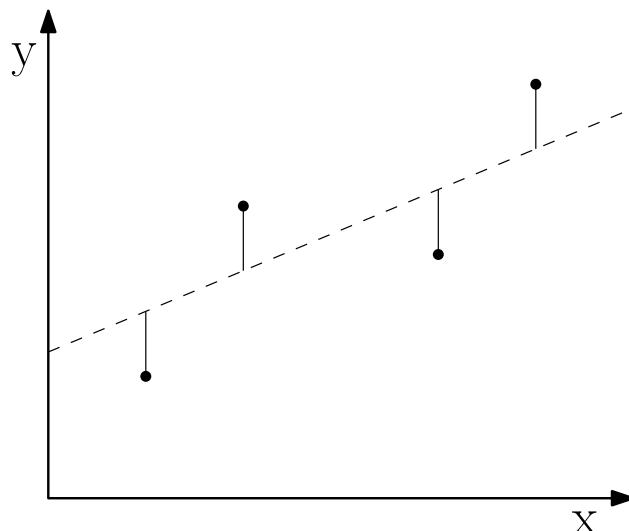
Blender: Image by Blender Foundation

Linear Fitting

$$\hat{y} = a + bx$$

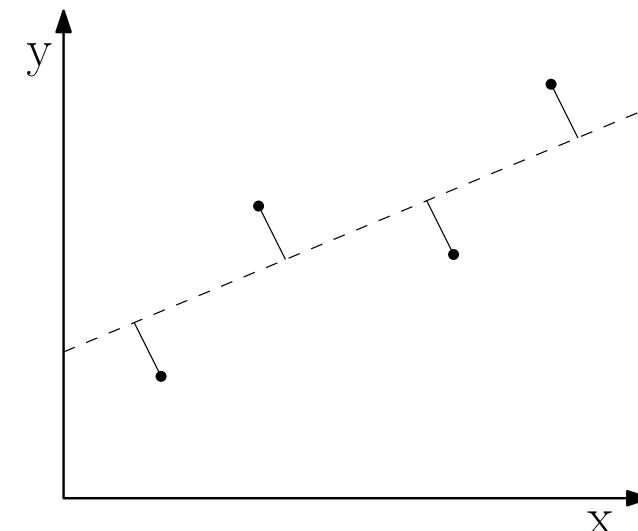
$$\text{minimize } R^2 = \sum_{i=1}^N r_i^2$$

vertical offsets



$$r_i = y_i - (a + bx_i)$$

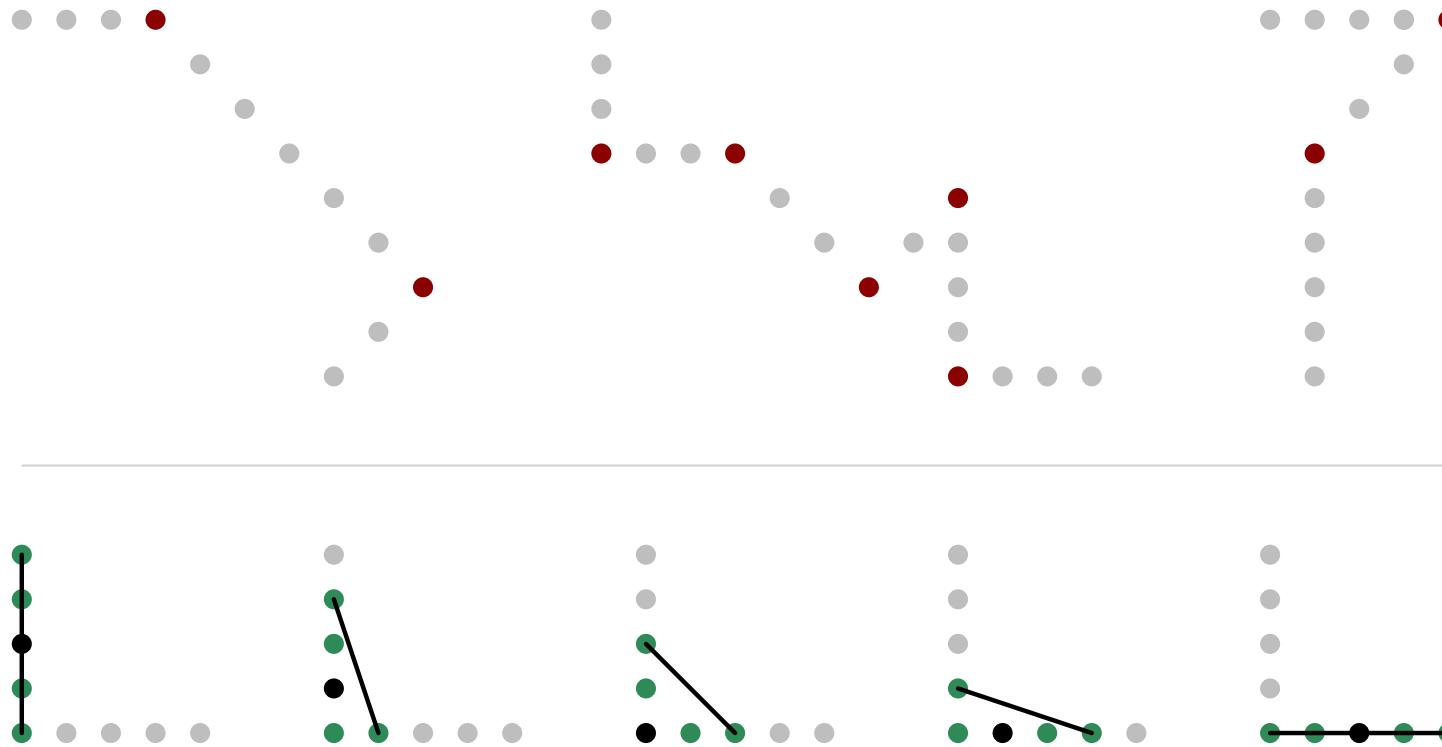
perpendicular offsets



$$r_i = \frac{|y_i - (a + bx_i)|}{\sqrt{1+b^2}}$$

Corner Detection

PDL, ShortStraw, Firefox's QuickGestures, etc.



$$G = s_1, \dots, s_n, \dots, s_N \mid s_n \in \{L, R, U, D\}$$

Graffiti & Unistrokes

A B C D E F G H I J K L M N O

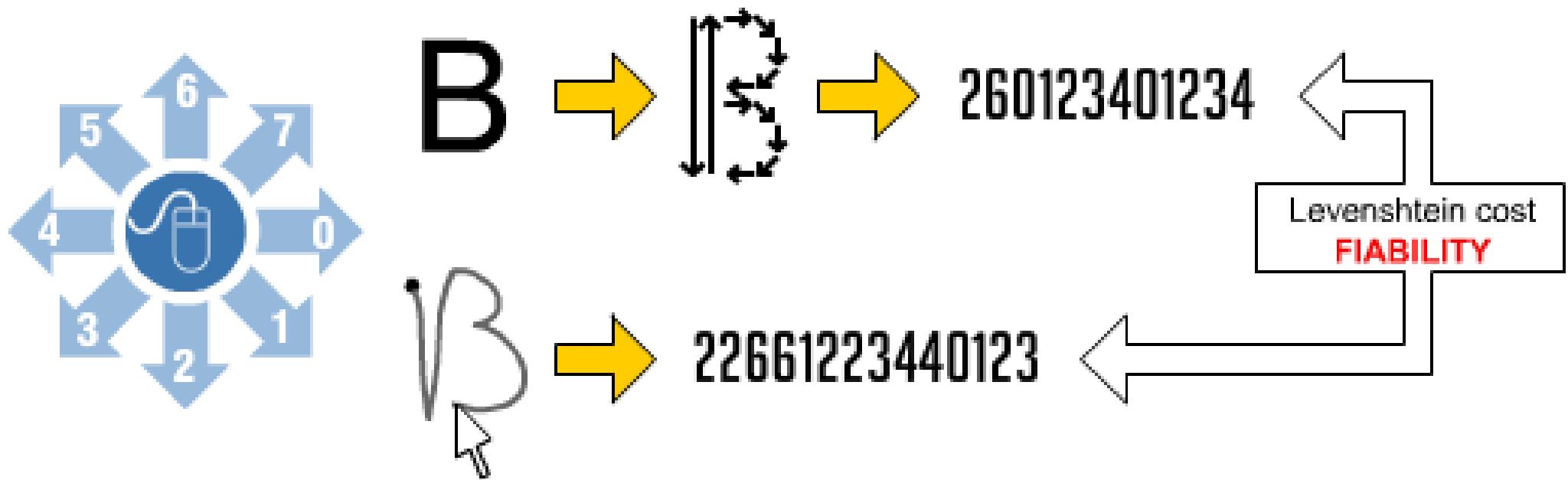
P Q R S T U V W X Y Z Space

A B C D E F G H I J K L M N O

P Q R S T U V W X Y Z (tap) Space

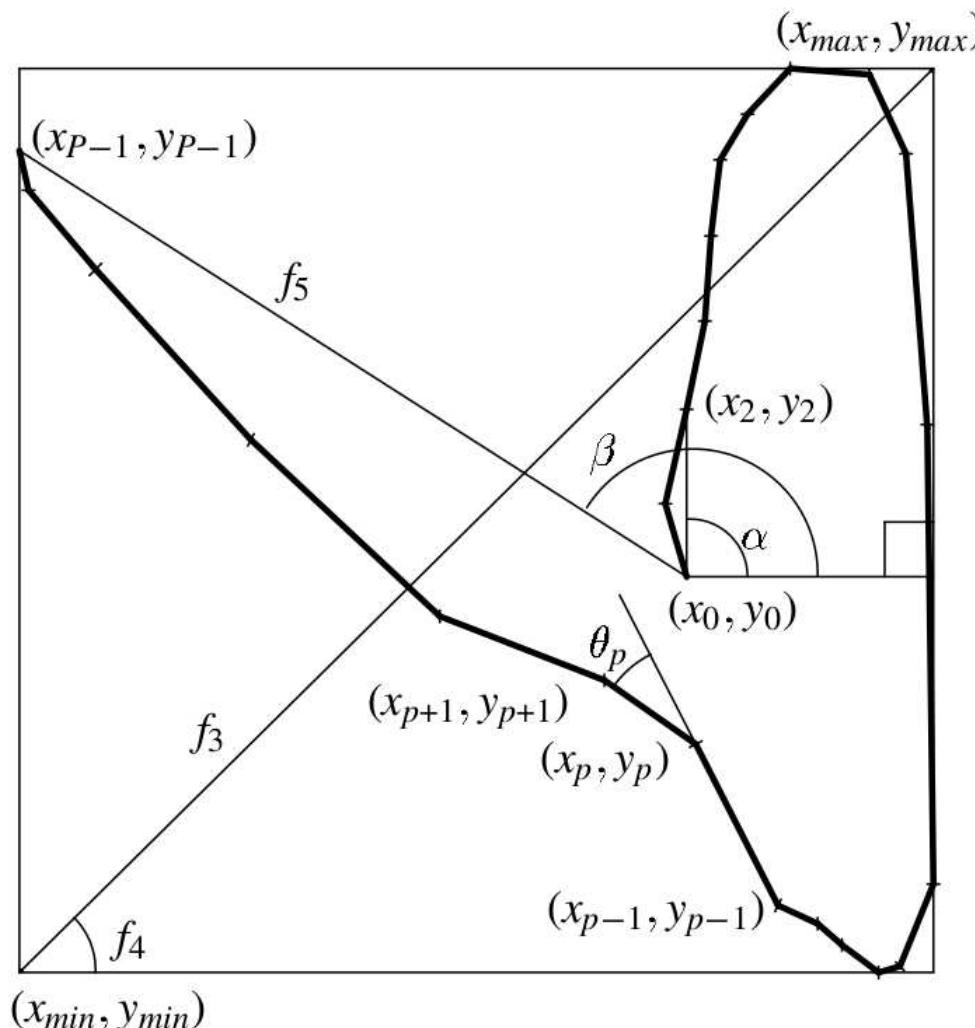
Comparison by [Castellucci and MacKenzie \(2008\)](#)

Graffiti & Unistrokes



Demo: <http://www.bytearray.org/wp-content/projects/mousegesture/GestureDemo.swf>

Rubine recognizer



by Rubine (1991)

Rubine recognizer

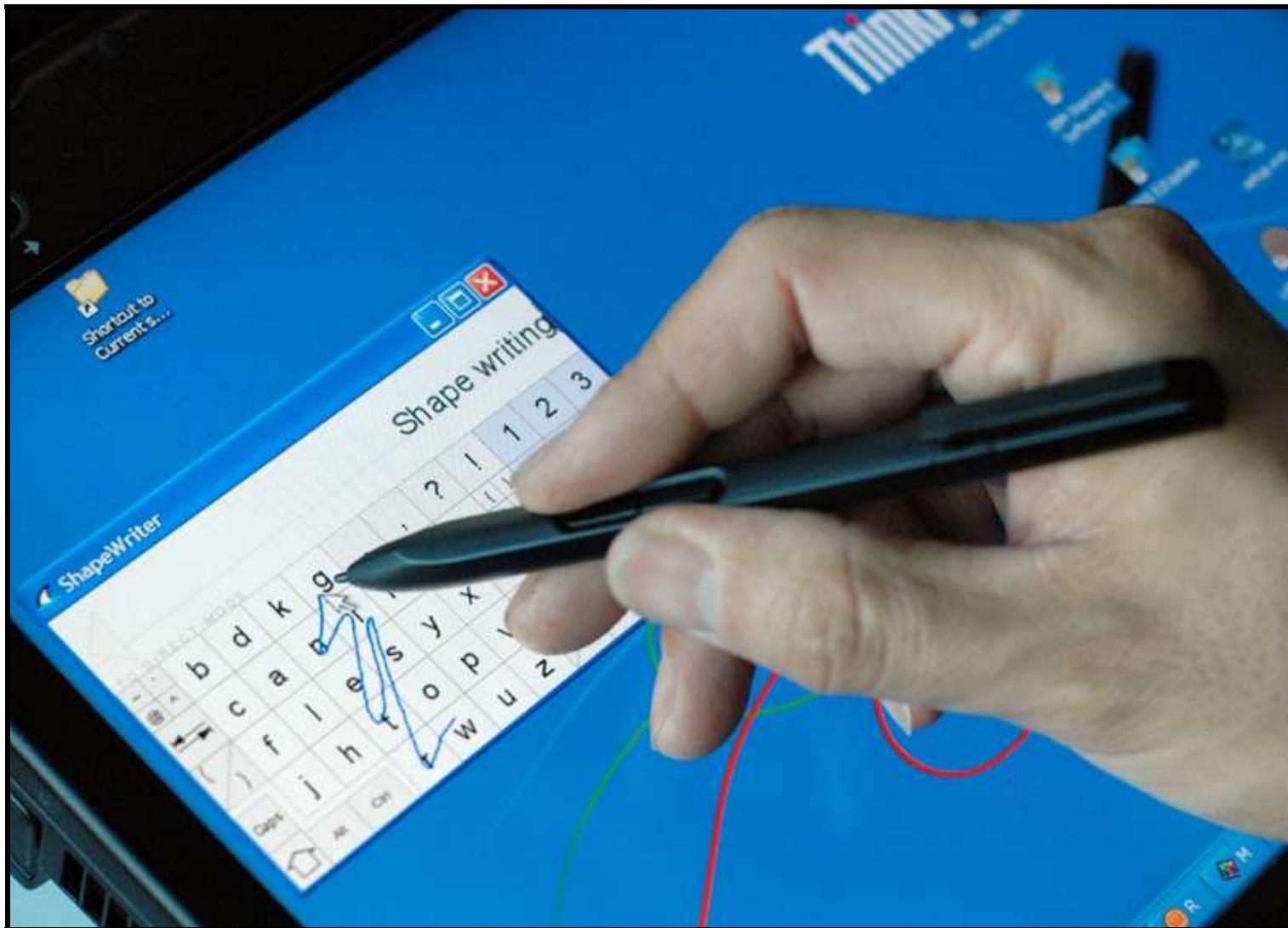
linear classifier using $F = 13$ stroke features

$$c = f(\mathbf{w}^T \mathbf{g}) = w_o + \sum_{i=1}^F \Sigma^{-1} \mu_i$$

$$w_0 = -\frac{1}{2} \sum_{i=1}^F w_i \mu_i$$

weight estimation:
perceptron, LSBF, LDA, SVM, logistic regression, etc.

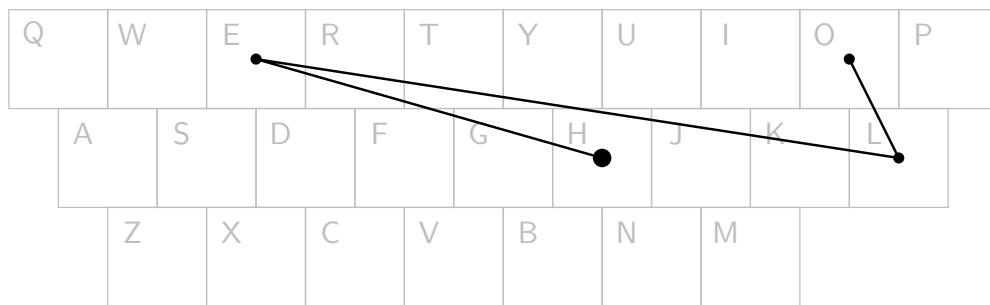
Shapewriting



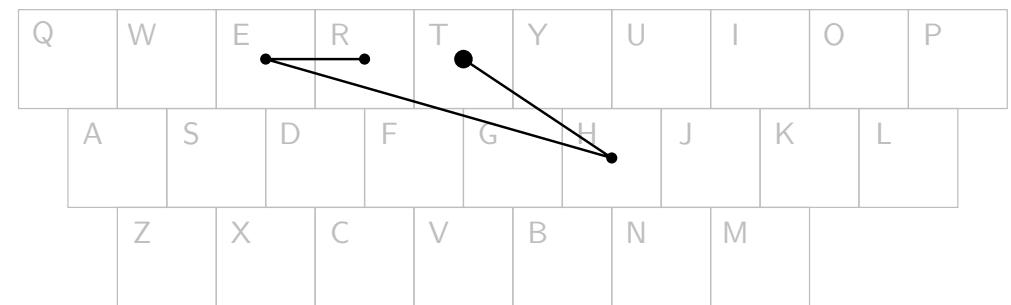
SHARK², by Kristensson and Zhai (2004)

Shapewriting

sokgraph of “hello”



sokgraph of “there”

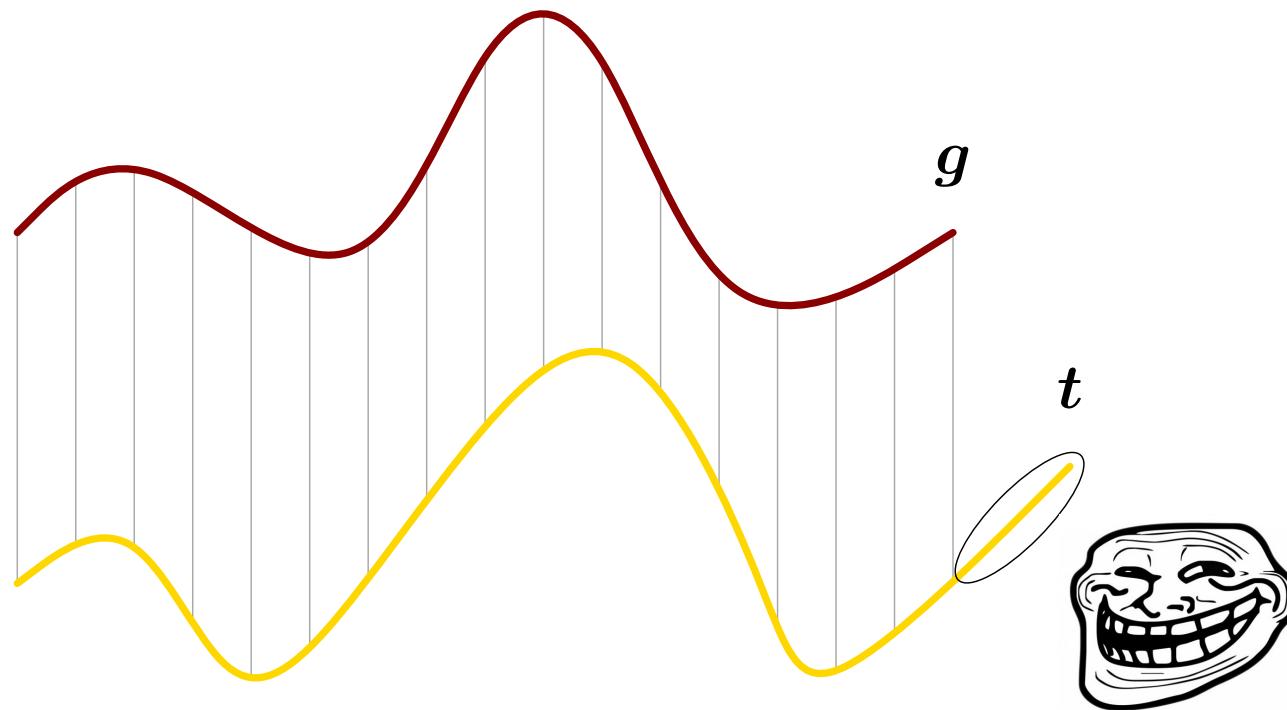


$$\hat{W} = \arg \max_W P(W|\mathbf{g})$$

$$\hat{W} = \arg \max_W \frac{P(\mathbf{g}|W)P(W)}{P(\mathbf{g})} = \arg \max_W P(\mathbf{g}|W)P(W)$$

Euclidean Matching

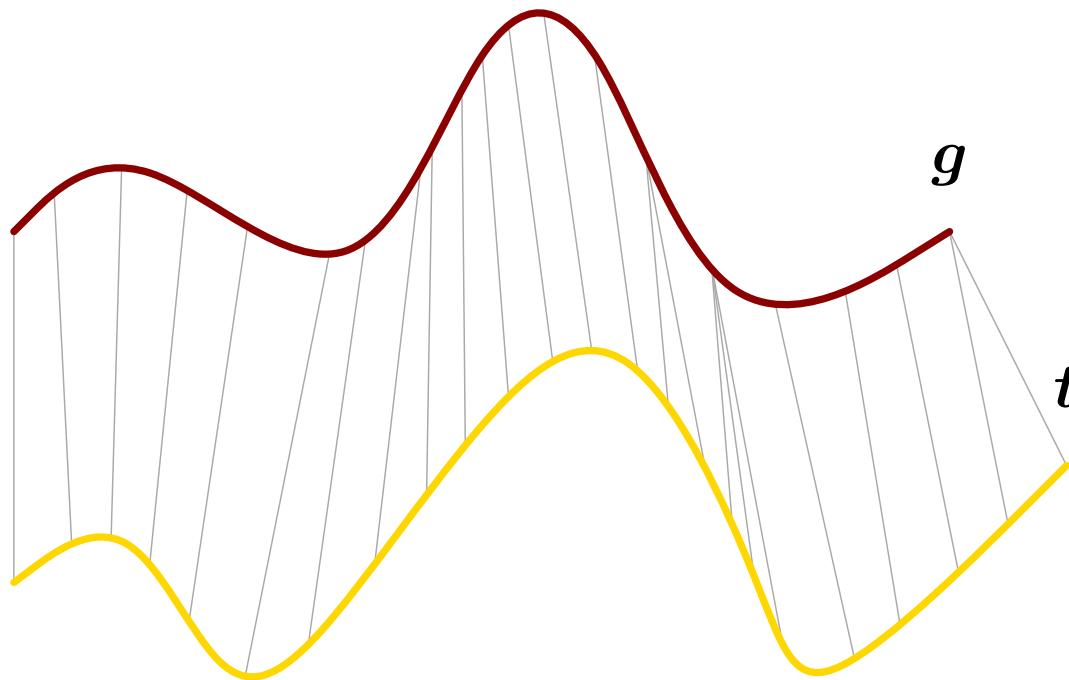
point-wise distances



$$D(\mathbf{g}, \mathbf{t}) = \frac{1}{|\mathbf{g}|} \sum_{i=1}^{|\mathbf{g}|} \|g_i - t_i\|$$

Elastic Matching

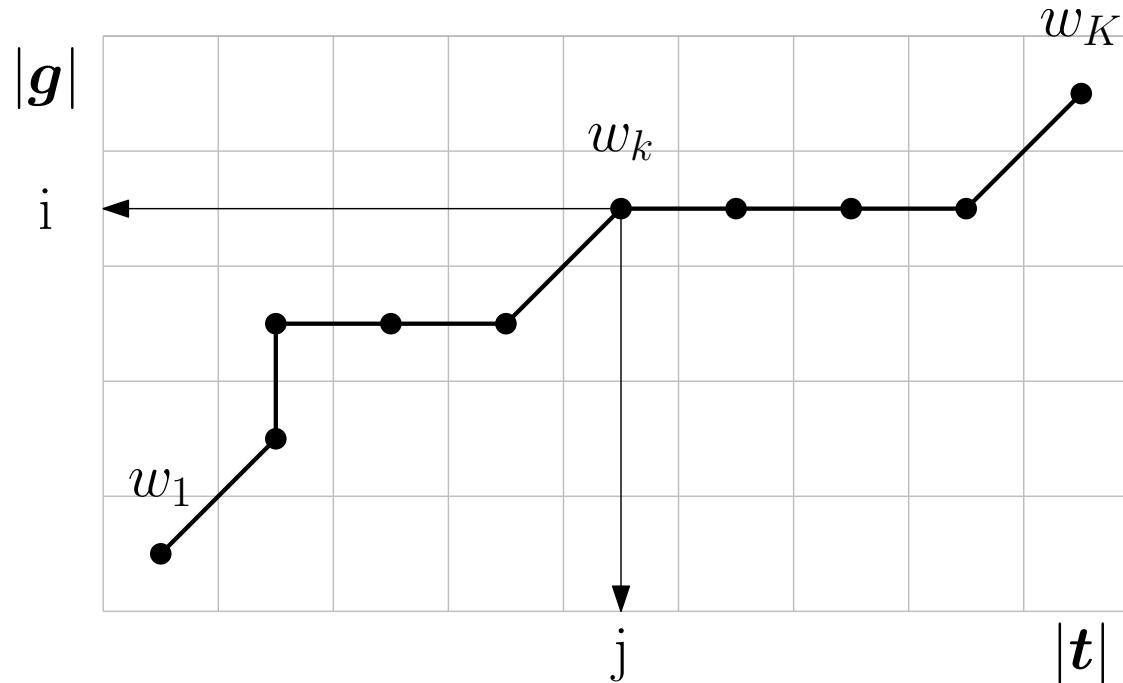
dynamic programming



$$D(g, t) = \min_{W \in \mathcal{W}} \frac{1}{|W|} \sum_{k=1}^{|W|} w_k$$

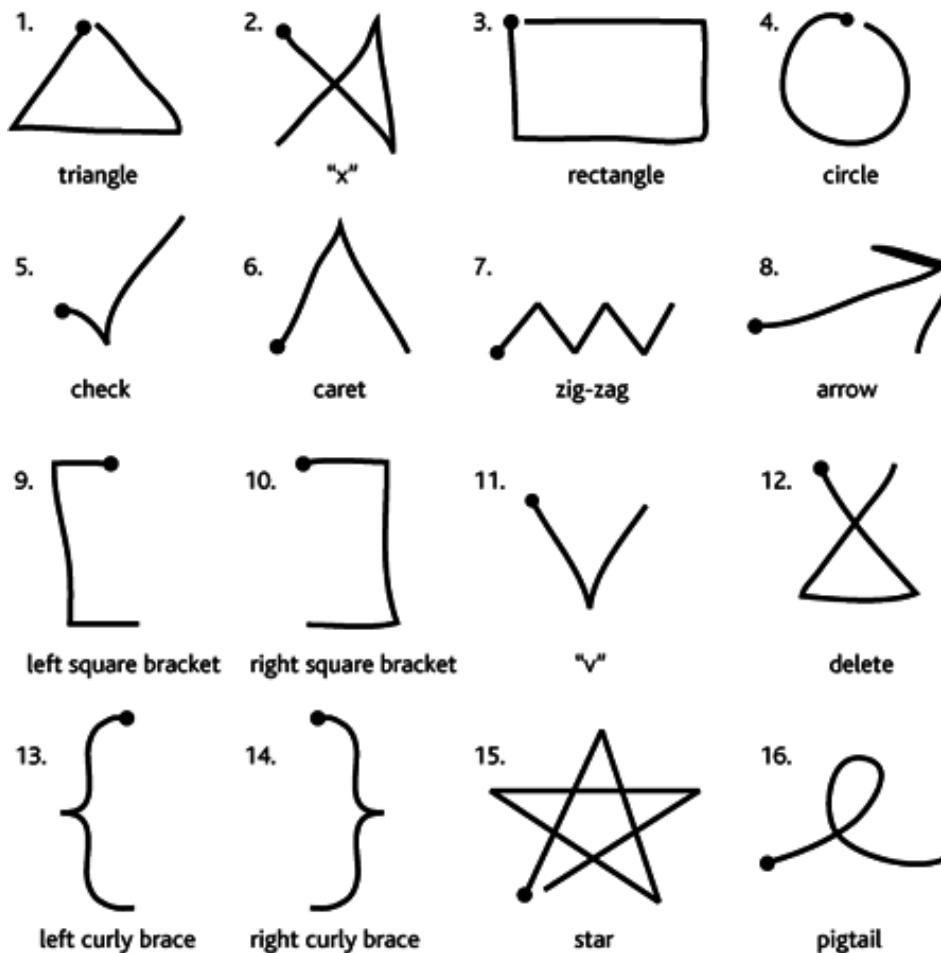
Elastic Matching

$$W = w_1, \dots, w_k, \dots, w_K$$



$$\gamma(i, j) = d(i, j) + \min\{\gamma(i - 1, j), \gamma(i, j - 1), \gamma(i - 1, j - 1)\}$$

The Dollar Family: \$1 recognizer

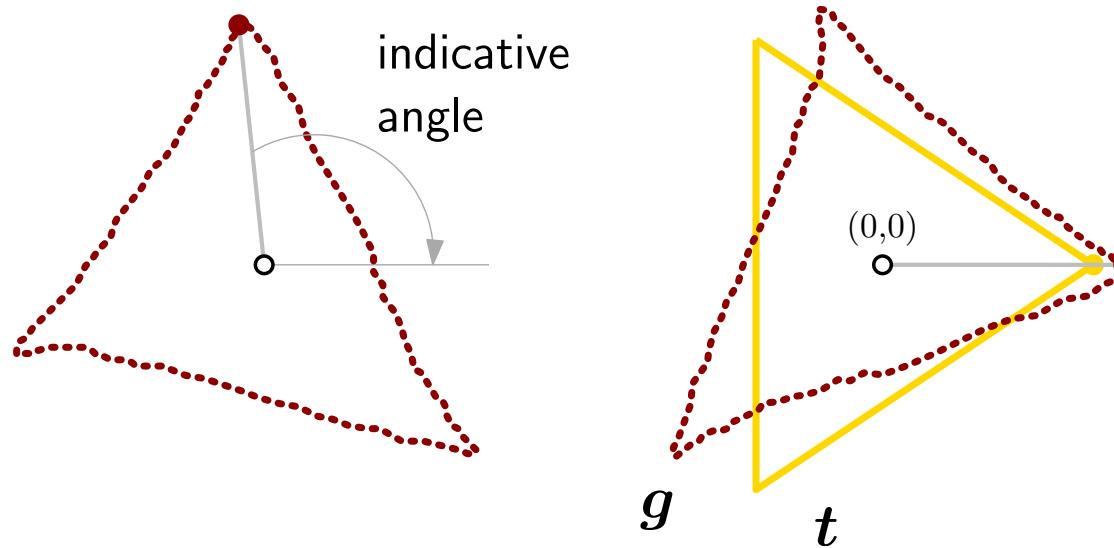


by Wobbrock et al. (2007)

<https://depts.washington.edu/aimgroup/proj/dollar/>

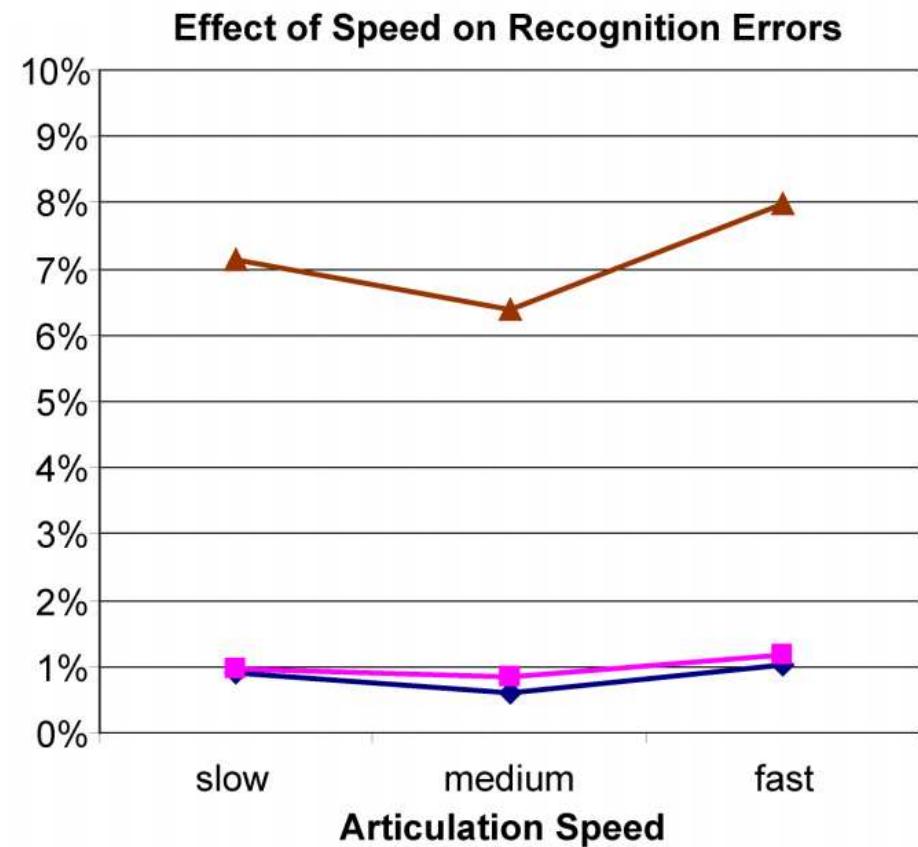
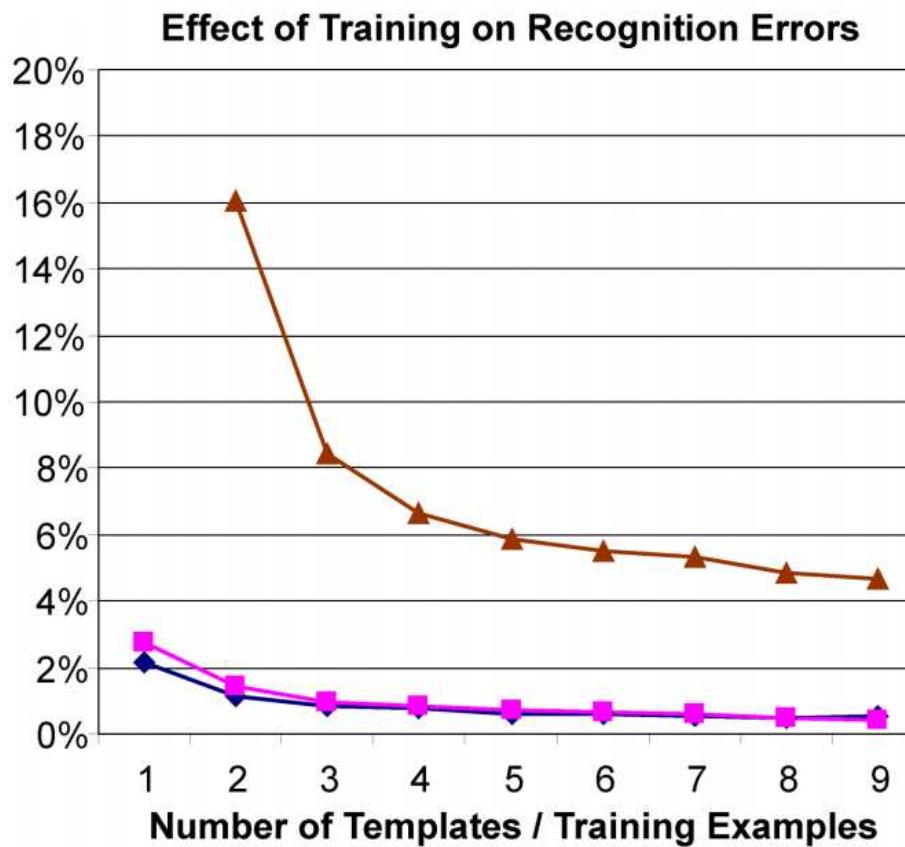
The Dollar Family: \$1 recognizer

Preprocessing: resampling, rotation, scaling, translation



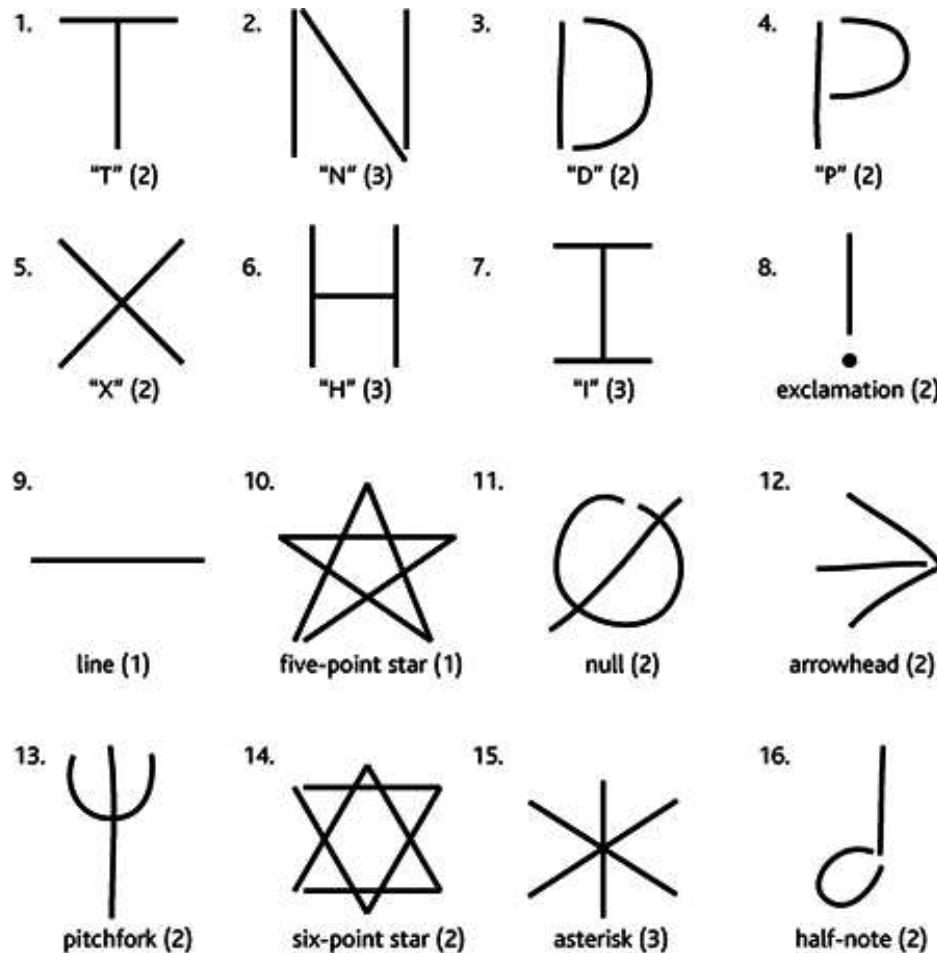
$$D(g, t) = \arg \min_{-\frac{\pi}{4} \leq \theta \leq \frac{\pi}{4}} \frac{1}{N} \sum_{i=1}^N \|g_i - t_i(\theta)\|$$

The Dollar Family: \$1 recognizer



► Rubine ■ \$1 ● DTW

The Dollar Family: \$N recognizer

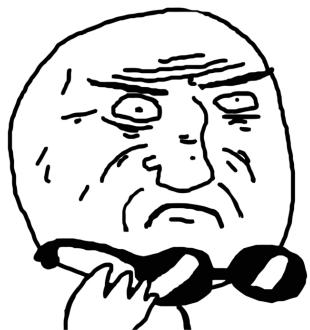
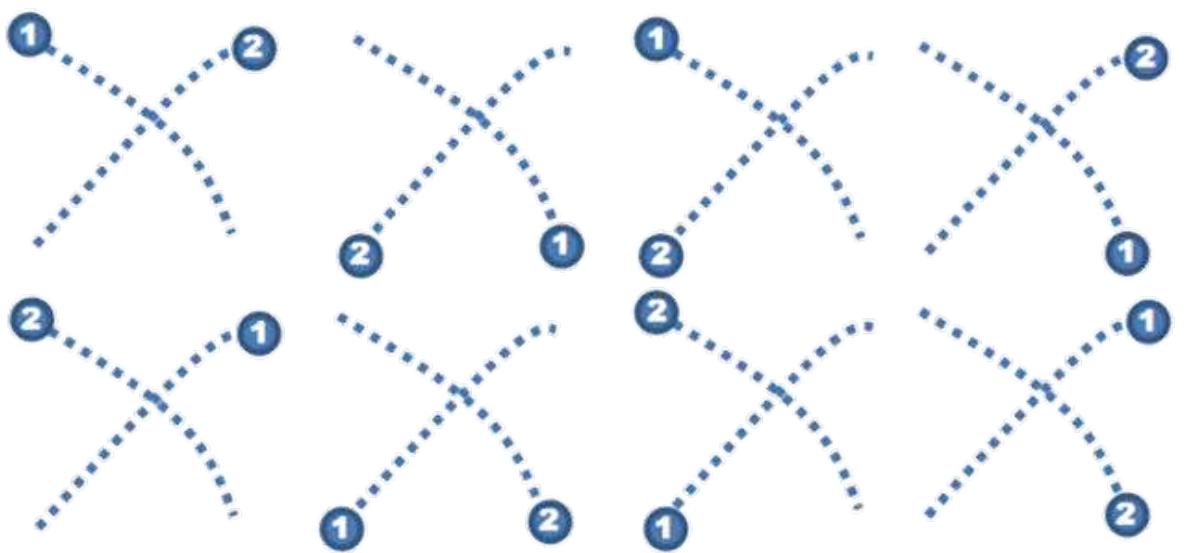


by Anthony and Wobbrock (2010)

<https://depts.washington.edu/aimgroup/proj/dollar/ndollar.html>

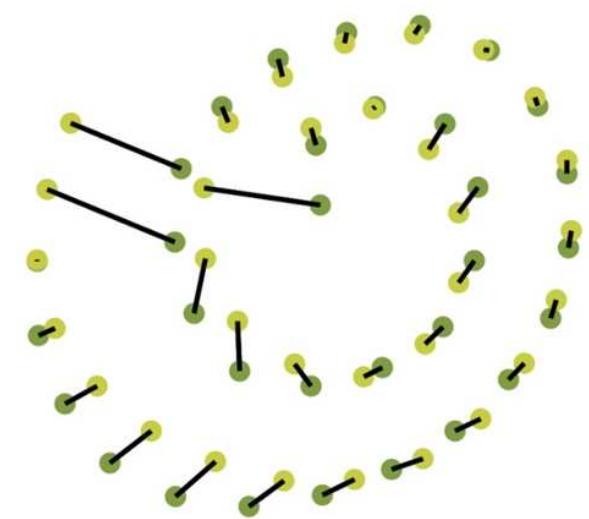
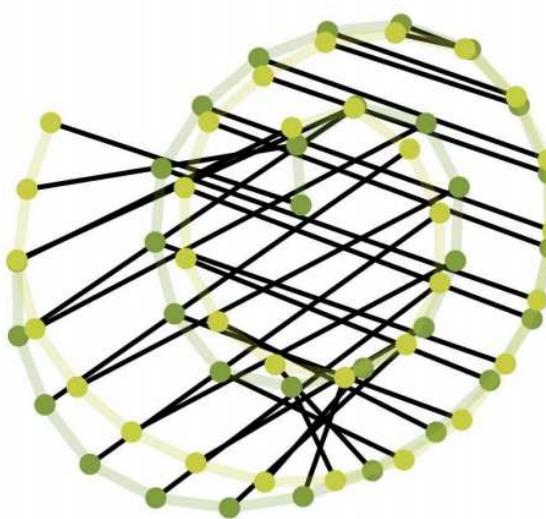
The Dollar Family: \$N recognizer

$\$N = \1 with combinatory overhead: $\mathcal{O}(n s 2^s)$ per template



Memory drained out with 20 templates ($n = 32$ pts)
in a quad-core computer with 4 GB RAM.

The Dollar Family: \$P recognizer



by Vatavu et al. (2012)

<https://depts.washington.edu/aimgroup/proj/dollar/pdollar.html>

The Dollar Family: \$P recognizer

variation of the Hungarian algorithm

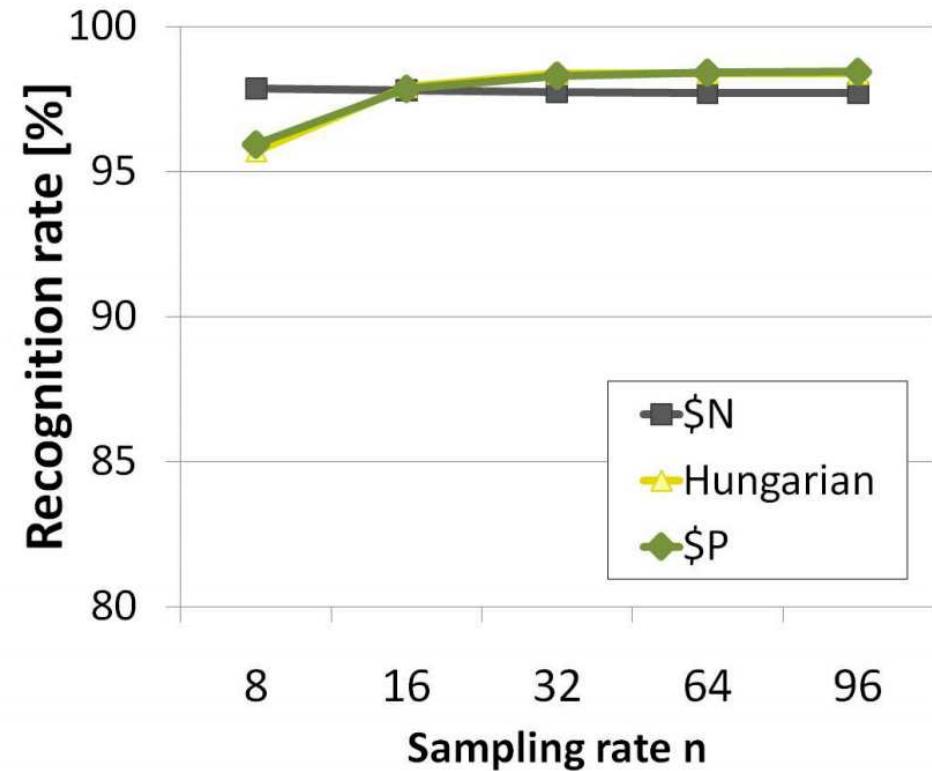
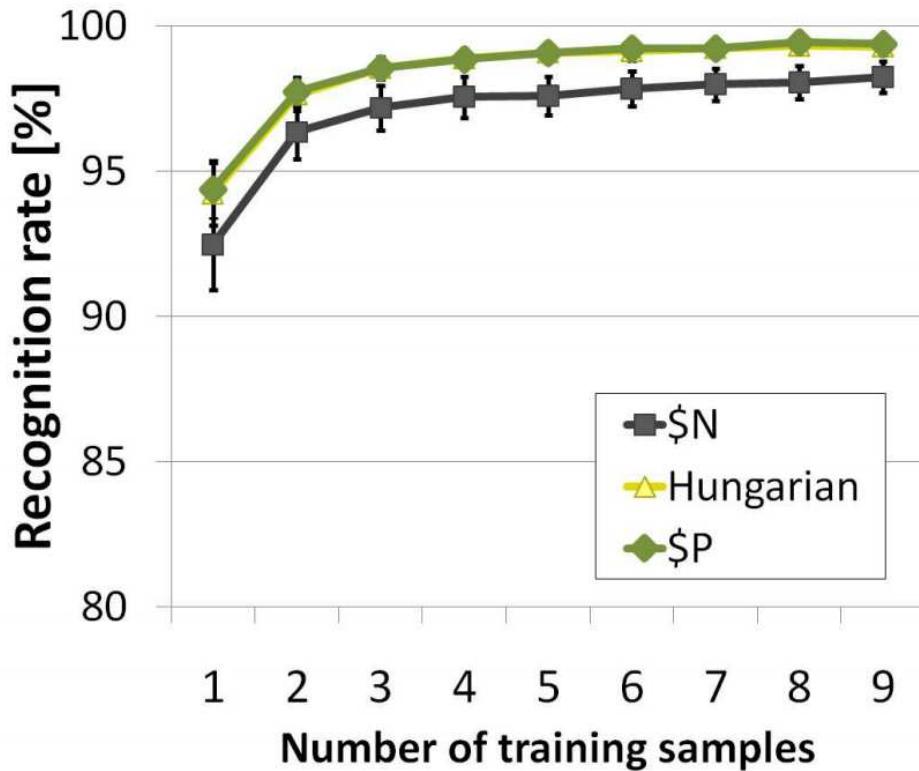
$$D(\mathbf{g}, \mathbf{t}) = \min \{\|\mathbf{g} - \mathbf{t}\|, \|\mathbf{t} - \mathbf{g}\|\}$$

Haussdorf's alternatives:

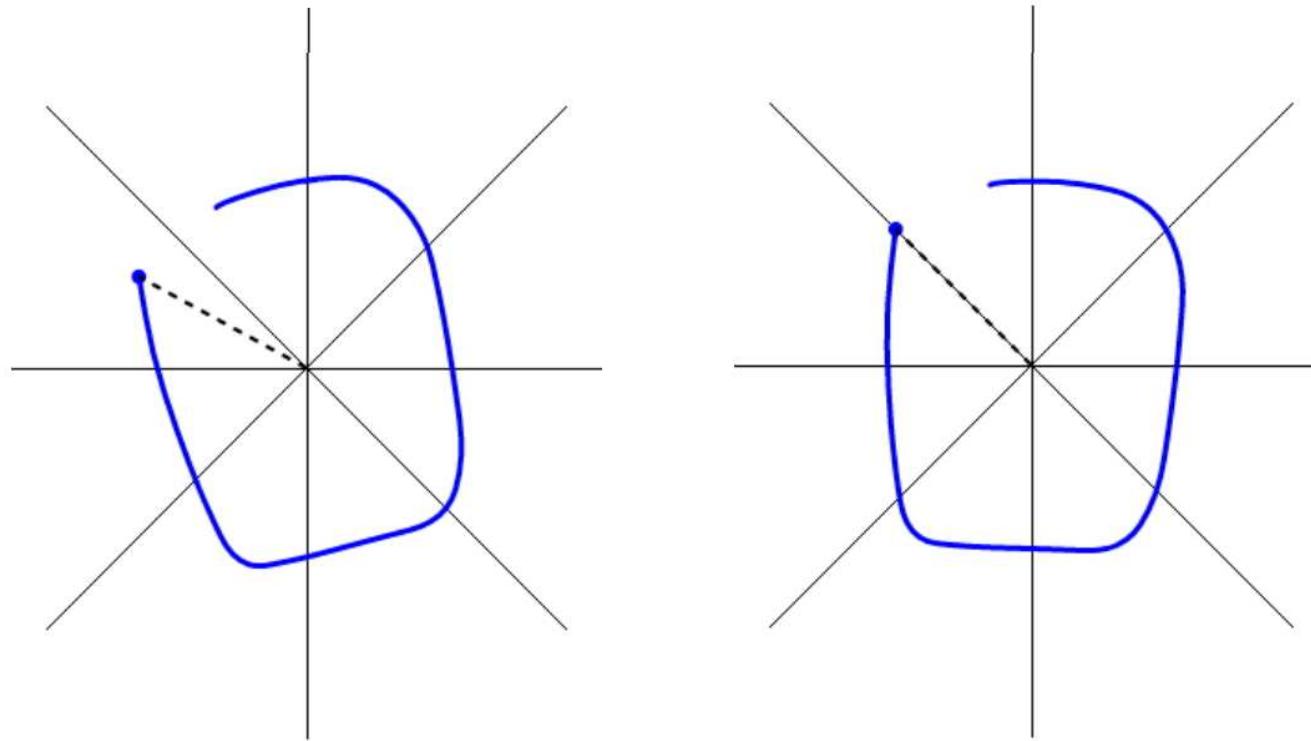
$$D(\mathbf{g}, \mathbf{t}) = \max_i \min_j \|\mathbf{g}_i - \mathbf{t}_j\|$$

$$D(\mathbf{g}, \mathbf{t}) = \frac{1}{N} \sum_{i=1}^N \min_j \|\mathbf{g}_i - \mathbf{t}_j\|$$

The Dollar Family: \$P recognizer



The Dollar Family: Protractor



by Li (2010)

The Dollar Family: Protractor

closed-form solution, minimum angular distance

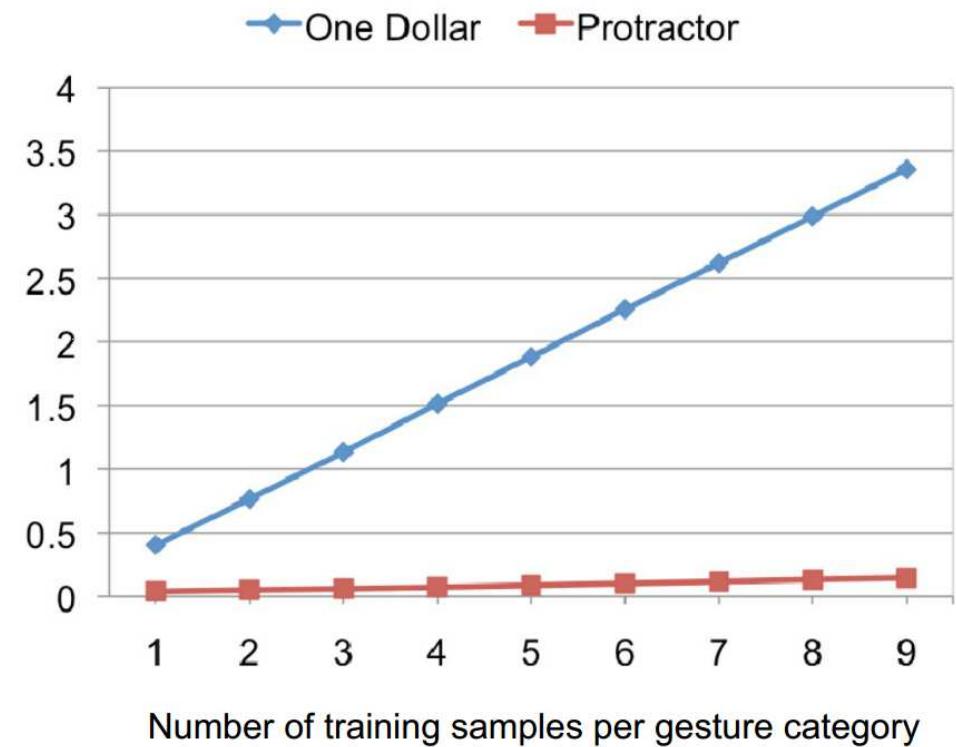
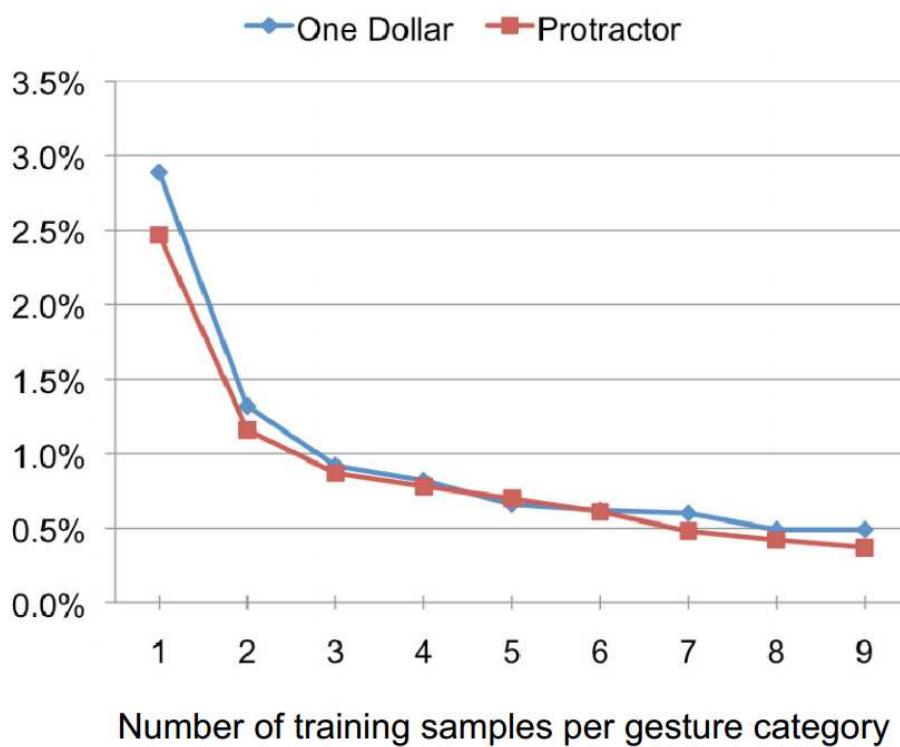
$$D(g, t) = \frac{1}{\arccos(a \cos \hat{\theta} + b \sin \hat{\theta})}$$

$$\hat{\theta} = \arctan \frac{b}{a}$$

$$a = \sum_{i=1}^N (x_{gi} x_{ti} + y_{gi} y_{ti})$$

$$b = \sum_{i=1}^N (x_{gi} y_{ti} - y_{gi} x_{ti})$$

The Dollar Family: Protractor



MinGestures for MIUIs

disambiguate gestures from text with high accuracy & performance

LABEL	ACTION	RESULT	LABEL	ACTION	RESULT
Substitute	Lorem Ipsu m <ins>a</ins> n	Lorem Ipsan	Split	Lorem	Lor em
Reject	Lorem Ipsum	Lorem ...	Validate	Lorem Ipsum	Lorem Ipsum
Merge	Lorem Ipsum	Lore <i>m</i> Ipsum	Undo	Lorem	Lorem Ipsum
Delete	Lorem Ipsum	Lorem	Redo	Lorem Ipsum	Lorem
Insert	Lorem Ipsum <ins>et</ins>	Lorem et Ipsum	Help	Lorem Ipsum	<help event>

by Leiva et al. (2014)

Demo: <http://cat.prhlt.upv.es/mg/>

MinGestures for MIUIs

disambiguating features:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N x_i y_i}$$

$$\Delta_x^- = \sum_{i=2}^N \max(x_{i-1} - x_i, 0)$$

$$\varphi = \frac{\max(\mathbf{x}) - \min(\mathbf{x})}{\max(\mathbf{y}) - \min(\mathbf{y})}$$

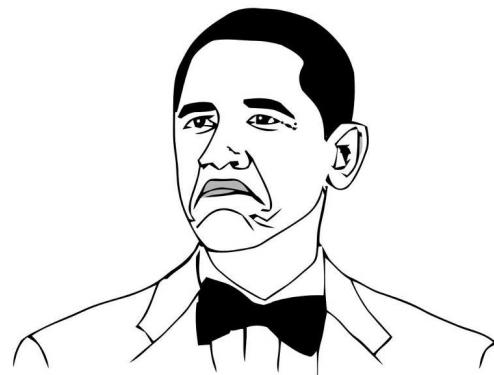
classification rule:

$$\theta = \frac{\hat{y} - b}{x} \pm \epsilon$$

MinGestures for MIUIs

System	E-pen		Mouse	
	training	test	training	test
\$1 recognizer	41.2	40.5	45.0	46.6
Marking Menus	18.5	19.5	12.8	13.1
Modified \$1	16.0	15.6	7.48	7.56
Rubine	14.6	14.1	15.4	15.7
MinGestures	0.77	1.32	0.26	0.43

Error rates in %



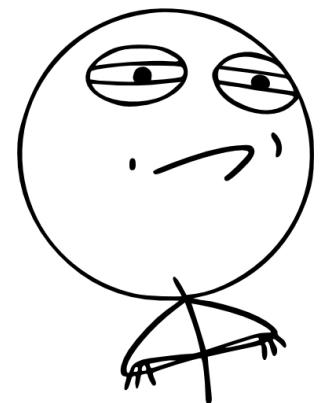
Recap

Takeaways

- Gestures aim to improve human-computer interaction
- Many recognition techniques suitable for different input devices
- Recognition trade-offs: design, setup, performance, accuracy
- Gestures should be **simple**:
 - For humans to *perform* and *recall*
 - For computers to *recognize*

TFM Proposals

1. Integration: free-form gestures in context
2. Error analysis & recovery:
When should the recognizer ask the user? How much to ask?
3. Segmentation: automatic gesture parts identification
4. Generation: grammar-based, kinematic theory, etc.





Bibliography

- L. ANTHONY AND J. O. WOBBROCK. A lightweight multistroke recognizer for user interface prototypes. In *Proc. GI*, 2010.
- O. BAU AND W. E. MACKAY. OctoPocus: A dynamic guide for learning gesture-based command sets. In *Proc. UIST*, 2008.
- S. J. CASTELLUCCI AND I. S. MACKENZIE. Graffiti vs. Unistrokes: An empirical comparison. In *Proc. CHI*, 2008.
- M. DAVIS AND T. ELLIS. The RAND tablet: A man-machine graphical communication device. In *Proc. AFIPS*, 1964.
- P. O. KRISTENSSON AND S. ZHAI. SHARK²: A large vocabulary shorthand writing system for pen-based computers. In *Proc. UIST*, 2004.
- G. P. KURTENBACH. *The design and evaluation of marking menus*. PhD thesis, University of Toronto, 1991.
- L. A. LEIVA, V. ALABAU, V. ROMERO, A. H. TOSELLI, AND E. VIDAL. Context-aware gestures for mixed-initiative text editing UIs. *Interacting with Computers*, 2014. To appear.
- Y. LI. Protractor: a fast and accurate gesture recognizer. In *Proc. CHI*, 2010.
- J. S. LIPSCOMB. A trainable gesture recognizer. *Pattern Recognition*, 24(9), 1991.

- R. PLAMONDON AND S. N. SRIHARI. On-line and off-line handwriting recognition: a comprehensive survey. *IEEE Transactions On Pattern Analysis And Machine Intelligence*, 22(1), 2000.
- D. H. RUBINE. *The Automatic Recognition of Gestures*. PhD thesis, Carnegie Mellon University, 1991.
- I. E. SUTHERLAND. Sketchpad: A man-machine graphical communication system. Tech. Report 296, Lincoln Laboratory, MIT, 1963.
- C. C. TAPPERT, C. Y. SUEN, AND T. WAKAHARA. The state of the art in on-line handwriting recognition. *IEEE Transactions On Pattern Analysis And Machine Intelligence*, 12(8), 1990.
- F. TIAN, F. LU, Y. JIANG, X. ZHANG, X. CAO, G. DAI, AND H. WANG. An exploration of pen tail gestures for interactions. *Int. J. Human-Computer Studies*, 71, 2012.
- R.-D. VATAVU, L. ANTHONY, AND J. O. WOBBROCK. Gestures as point clouds: a \$P recognizer for user interface prototypes. In *Proc. ICMI*, 2012.
- J. O. WOBBROCK, A. D. WILSON, AND Y. LI. Gestures without libraries, toolkits or training: A \$1 recognizer for user interface prototypes. In *Proc. UIST*, 2007.
- S. ZHAI, P. O. KRISTENSSON, C. APPERT, T. H. ANDERSEN, AND X. CAO. Foundations and trends in Human–Computer Interaction. *Foundational Issues in Touch-Surface Stroke Gesture Design – An Integrative Review*, 5(2), 2012.

Videography

\$1 GESTURE RECOGNITION. <http://www.youtube.com/watch?v=ept2Z1UVxfw>.
HUMANTENNA. <http://www.youtube.com/watch?v=7lRnm2oFGdc>.
LEAP MOTION. http://www.youtube.com/watch?v=_d6KuiuteIA.
MARKING MENUS. <http://www.youtube.com/watch?v=dtH9GdFSQaw>.
MYO. <http://www.youtube.com/watch?v=oWu9TFJjHaM>.
OCTOPOCUS. <http://vimeo.com/2116172>.
PEN TAIL. <http://iel.iscas.ac.cn/~fengt/videos/pentailgestures.wmv>.
RAND TABLET. <http://www.youtube.com/watch?v=LLRy4Ao62ls>.
SHAPEWRITING. <http://www.youtube.com/watch?v=WtlyuuYmFNO>.
SHUTERLAND'S SKETCHPAD. http://www.youtube.com/watch?v=USyoT_Ha_bA.
SKINPUT. <http://www.youtube.com/watch?v=g3XPUDW9Ryg>.
WACOM GESTURES. <http://www.youtube.com/watch?v=hKQFqEVK81M>.