

LB_Keogh Supports Exact Indexing of Shapes under Rotation Invariance with Arbitrary Representations and Distance Measures

Keogh, Wei, Xi, Lee & Vlachos

Come, we shall learn of the indexing of shapes

Set forth these figures as I have conceived their shape...*

Outline of Talk

- The utility of shape matching
- Shape representations
- Shape distance measures
- Lower bounding rotation invariant measures with the LB_Keogh
- Accuracy experiments
- Efficiency experiments
- Conclusions

Paradise Comix XVII 85

The Utility of Shape Matching I

...discovering insect mimicry, clustering petroglyphs, finding unusual arrowheads, tracking fish migration, finding anomalous fruit fly wings...

The Utility of Shape Matching II

...automatically annotating old manuscripts, mining medical images, biometrics, spatial mining of horned lizards, indexing nematodes...

Shape Representations I

For virtually all shape matching problems, **rotation** is the problem

If I asked you to group these reptile skulls, **rotation** would not confuse you

There are two ways to be rotation invariant

- 1) Landmarking: Find the one "true" rotation
- 2) Rotation invariant features

Landmarking

- **Domain Specific Landmarking**
Find some fixed point in your domain, eg. the nose on a face, the stem of leaf, the tail of a fish ...
- **Generic Landmarking**
Find the major axis of the shape and use that as the canonical alignment

Best Rotation Alignment

Generic Landmark Alignment

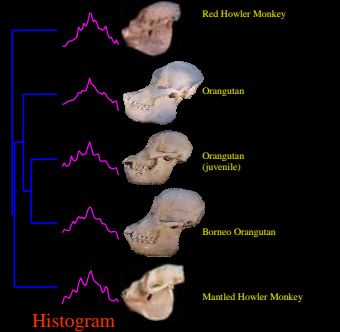
The only problem with landmarking is that it does not work

Rotation invariant features

Possibilities include:
Ratio of perimeter to area, fractal measures, elongatedness, circularity, min/max/mean curvature, entropy, perimeter of convex hull and histograms



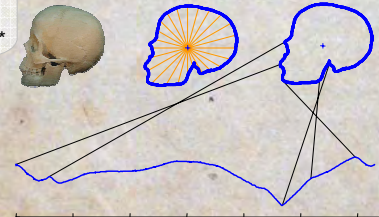
The only problem with rotation invariant features is that in throwing away rotation information, you must invariably throw away useful information



We can convert shapes into a 1D signal. Thus can we remove information about *scale* and *offset*.

Rotation we must deal with in our algorithms...

...so it seemed to change its shape, from running lengthwise to revolving round...*



There are many other 1D representations of shape, and our algorithm can work with any of them

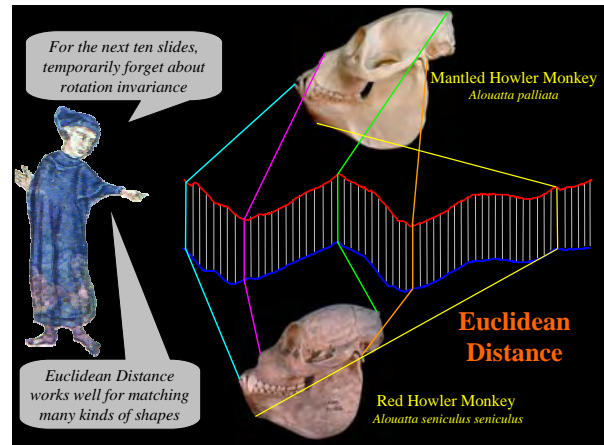
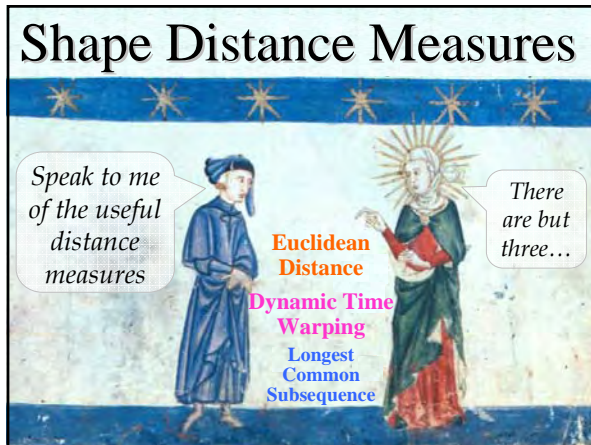
*Dante Alighieri, The Divine Comedy Paradise - Canto XXX, 90.

Shape Distance Measures

Speak to me of the useful distance measures

Euclidean Distance
Dynamic Time Warping
Longest Common Subsequence

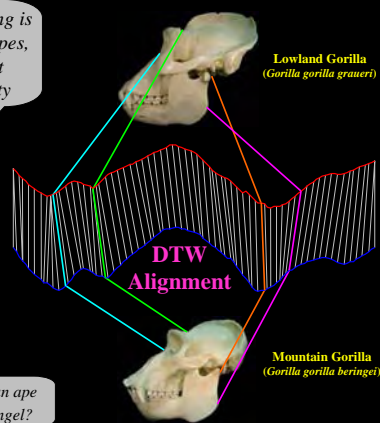
There are but three...



Dynamic Time Warping is useful for natural shapes, which often exhibit intraclass variability



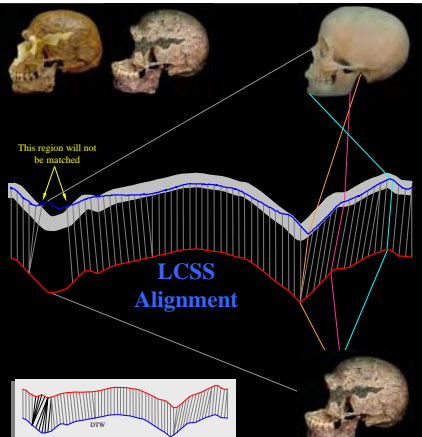
Is man an ape or an angel?



Matching skulls is an important problem



LCSS can deal with missing or occluded parts



For brevity, we will only give details of Euclidean distance in this talk

However, the main point of our paper is that the same idea works for DTW and LCSS with no overhead

We will present empirical results that do show that DTW can be significantly better than Euclidean distance

Euclidean Distance Metric

Given two time series $Q = q_1 \dots q_n$ and $C = c_1 \dots c_n$, the Euclidean distance between them is defined as:

$$D(Q, C) = \sqrt{\sum_{i=1}^n (q_i - c_i)^2}$$

I notice that you Z-normalized the time series first

The next slide shows a useful optimization

Early Abandon Euclidean Distance

During the computation, if current sum of the squared differences between each pair of corresponding data points exceeds r^2 , we can safely **abandon** the calculation

I see, because incremental value is always a lower bound to the final value, once it is greater than the best-so-far, we may as well abandon

Abandon all hope ye who enter here

Most indexing techniques work by grouping objects into logical units, and defining a lower bound distance to the units

Here we will use "wedges" as the logical unit, and LB_Keogh as the lower bound distance

For example, for indexing cities we can use MBRs and the classic MIN-DIST function of Guttman

Wedge

Suppose two shapes get converted to time series...

Having candidate sequences C_1, \dots, C_k , we can form two new sequences U and L :

$$U_i = \max(C_{1i}, \dots, C_{ki})$$

$$L_i = \min(C_{1i}, \dots, C_{ki})$$

They form the smallest possible bounding envelope that encloses sequences C_1, \dots, C_k .

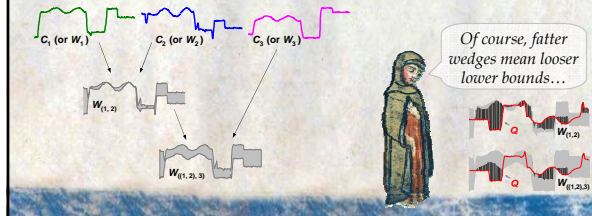
We call the combination of U and L a wedge, and denote a wedge as $W = \{U, L\}$

A lower bounding measure between an arbitrary query Q and the set of candidate sequences contained in a wedge W , is the LB_Keogh

$$LB_Keogh(Q, W) = \begin{cases} (q_i - U_i)^2 & \text{if } q_i > U_i \\ (q_i - L_i)^2 & \text{if } q_i < L_i \\ 0 & \text{otherwise} \end{cases}$$

Generalized Wedge

- Use $W_{(1,2)}$ to denote that a wedge is built from sequences C_1 and C_2 .
- Wedges can be hierarchically nested. For example, $W_{((1,2),3)}$ consists of $W_{(1,2)}$ and C_3 .

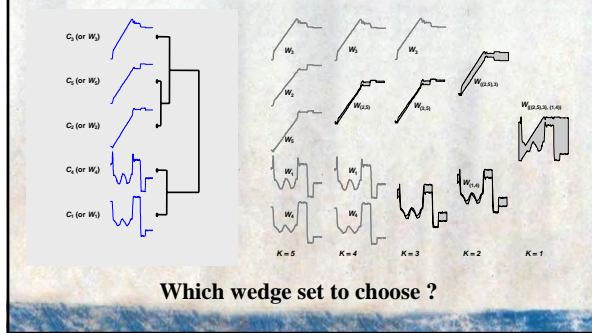


We are finally ready to explain our idea for rotation invariance, an idea we have sidestepped to this point. Suppose we have a shape as before...

We can create every possible rotation of the shape, by considering every possible circular shift of the time series, as shown at my left... But we already know how to index such time series by using wedges! We just need to figure out the best wedge making policy.

It sucks being a grad student

Hierarchical Clustering



Once we have all possible rotations of all the objects we want to index inserted into wedges, we can simply use any LB_Keogh indexer

Since the introduction of LB_Keogh indexing at this conference 4 years ago, at least 50 groups around the world have used/extended/adapted the idea, making this work easily reimplementable

What are the disadvantages of using LB_Keogh?

There are Nun

"LB_Keogh has provided a convincing lower bound" T. Rath
 "LB_Keogh can significantly speed up DTW." Suzuki
 "LB_Keogh is the best..." Zhou & Wong
 "LB_Keogh offers the tightest lower bounds". M. Cardle.
 "LB_Keogh makes retrieval of time-warped time series feasible even for large data sets". Muller et. al.
 "LB_Keogh can be effectively used, resulting in considerably less number of DTW computations." Karydis
 "exploiting LB_Keogh, we can guarantee indexability". Bartolini et. al.
 "LB_Keogh, the best method to lower bound.." Capitani.
 "LB_Keogh is fast, because it cleverly exploits global constraints that appear in dynamic programming" Christos Faloutsos.

By using the LB_Keogh framework, we can leverage off the wealth of work in the literature

All our Experiments are Reproducible!

People that do irreproducible experiments should be boiled alive

Agreed! All our data is publicly available

www.cs.uci.edu/~eamonn/shape/

We tested on many diverse datasets

...and I recognized the face [¶]

Leaf of mine, in whom I found pleasure [†]

...as a fish dives through water [£]

...the shape of that cold animal which stings and lashes people with its tail ^{*}

†Purgatorio - Canto XXIII, ¶Paradiso - Canto XXV, £L'Inferno - Canto XXIV, *L'Inferno - Canto XXIV

Name	Classes	Instances	Euclidean Error (%)	DTW Error (%) (R)	Other Techniques
Face	16	2240	3.839	3.170 (3)	
Swedish Leaves	15	1125	13.33	10.84 (2)	
Chicken	5	446	19.96	19.96(1)	20.5 Discrete strings
MixedBag	9	160	4.375	4.375(1)	Chamfer 6.0, Hausdorff 7.0
OSU Leaves	6	442	33.71	15.61 (2)	
Diatoms	37	781	27.53	27.53(1)	26.0 Morphological Curvature Scale Spaces
Plane	7	210	0.95	0.0 (3)	0.55 Markov Descriptor
Fish	7	350	11.43	9.71 (1)	36.0 Fourier /Power Cepstrum

Note that DTW is sometimes worth the little extra effort

Implementation details should not matter, for example the results reported should be the same if reimplemented in Ret Hat Linux

We therefore use a cost model that is independent of hardware/software/buffer size etc. See the paper for details

We compare to brute force, and were possible a Fourier based approach (it can't handle DTW)

Main Memory Experiments

- Projectile point database
- Increasingly larger datasets
- One-nearest-neighbor queries

Euclidean

DTW

Number of objects in database (n)

Indexing Experiments

- Projectile point/Heterogenous databases
- Increasingly large dimensionality
- One-nearest-neighbor queries

Projectile Points

Heterogeneous

Fraction of objects retrieved

Dimensionality

Wedge: Euclidean

Wedge: DTW

... from its stock this tree was cultivated *

All these are in the genus *Cercopithecus*, except for the skull identified as being either a Vervet or Green monkey, both of which belong in the Genus of *Chlorocebus* which is in the same Tribe (*Cercopithecoidea*) as *Cercopithecus*.

Tribe *Cercopithecoidea*

Cercopithecus

En Brown Monkey, *Cercopithecus neglectus*, Moustached Monkey, *Cercopithecus cephus*, Red-tailed Monkey, *Cercopithecus ascanus*, *Chlorocebus*

Green Monkey, *Chlorocebus sabaeus*, Vervet Monkey, *Chlorocebus pygmythrix*

These are the same species

Homopithecus hoolock (Hoolock Gibbon)

These are in the Genus *Pongo*

All these are in the family *Cebidae*

Family *Cebidae* (New World monkeys)

Subfamily *Aotinae*

Aotus trivirgatus

Subfamily *Pitheciinae* ailla

Black Howler Saki, *Chiropotes satanas*, White-headed Saki, *Chiropotes albinatus*

All these are in the tribe

Tribe *Papionini*

Genus *Papio* - baboons

Genus *Mandrillus* - Mandrill

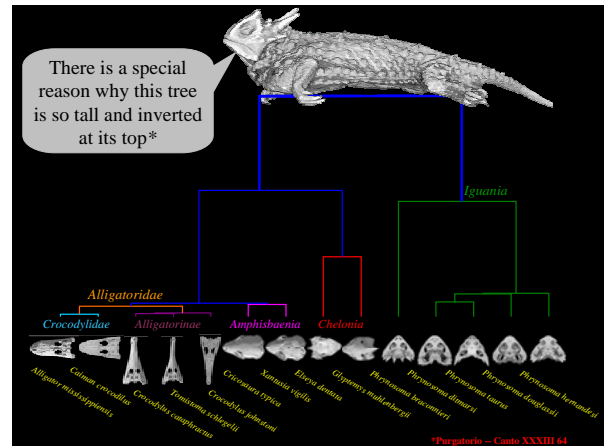
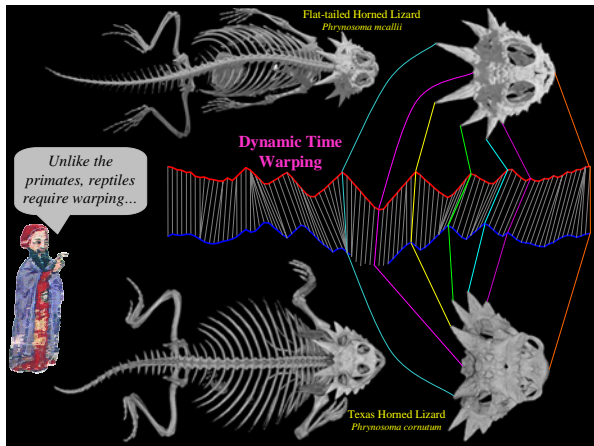
These are in the family *Lemniscata*

These are in the genus *Alouatta*

These are in the same species

Homo sapiens (Humans)

*Purgatorio - Canto XXIV 117



Petroglyph Mining

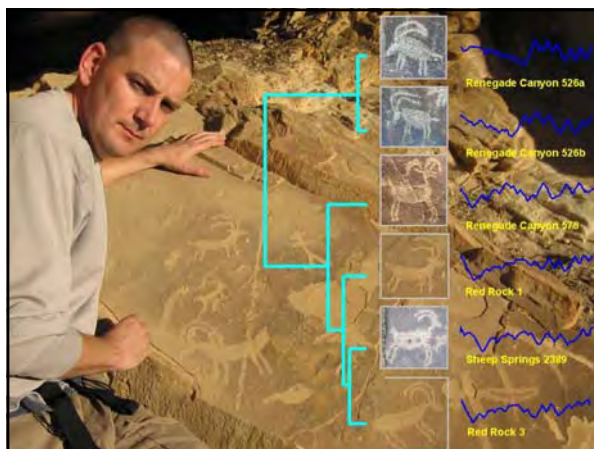
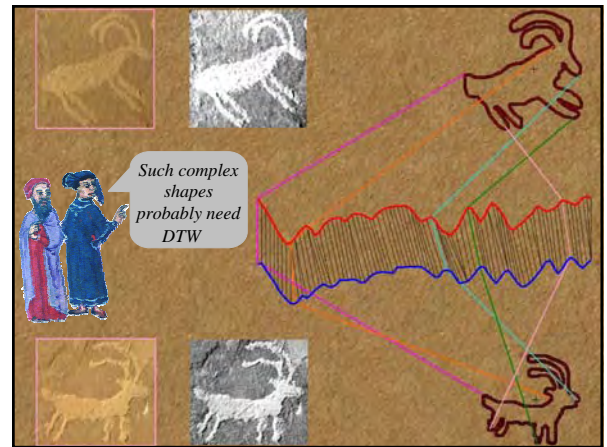
Petroglyphs are images incised in rock, usually by prehistoric peoples. They were an important form of pre-writing symbols, used in communication from approximately 10,000 B.C.E. to modern times. [Wikipedia](#)

- They appear worldwide
- Over a million in America alone
- Surprisingly little known about them

who so sketched out the shapes there?*

.. they would strike the subtlest minds with awe*

*Purgatorio -- Canto XII 6



Future Work: Data Mining

We did not want to work on shape data mining until we could do fast matching, that would have been ass backwards

.. so similar in act and coloration that I will put them both to one*

*Inferno -- Canto XXIII 29

Phylogenetic tree showing relationships between *Limenitis* and *Danaus* species.

Questions?

Feel free to email us with questions
Eamonn Keogh: Project Leader
eamonn@cs.ucr.edu

Li Wei: Lower Bounding
lw@cs.ucr.edu

Michail Vlachos: Public Nudity and Index Structures
vlachos@us.ibm.com

Sang Hee Lee: Anthropology and Primatology
shlee@ucr.edu

Xiaopeng Xi: Image Processing
xxi@cs.ucr.edu