

# Multidimensional Compression With Pattern Matching

### Introduction

Scientific data is usually collected to a greater degree of precision than is significant. Can we take advantage of this feature to reduce storage? We are working on a dictionary-based compression algorithm that finds statistically similar 1-dimensional data blocks. Here, we consider multidimensional similarity based compression methods as measured by peak signal-to-noise ratios (PSNR) and runtime over different compression levels.

## **Similarity Measures**

• Data points within partitioned time series are assumed to have the statistical property of exchangeability

•Kolmogorov-Smirnov (KS) statistical similarity test can then be performed



Figure 1: Schematic of a dictionary-based data compression algorithm known as IDEALEM (Implementation of Dynamic Extensible Adaptive Locally Exchangeable Measures)

- The KS test is not multidimensional
- •We propose alternative similarity measures: **Dynamic Time Warp** (DTW) and Minimum Jump Cost (MJC) • Consider two time series **x** and **y**

 $\mathbf{x} = (4, 4, 3, 2, 4, 0)$  $\mathbf{y} = (7, 5, 6, 4, 3, 4)$ 







```
A(\mathbf{x},\mathbf{y}) =
```



Figure 3: Visualization of a DTW distance measurement









Olivia Del Guercio<sup>1</sup>, Rafael Orozco<sup>2</sup>, Alexander Sim<sup>3</sup>, John Wu<sup>3</sup> <sup>1</sup>Scripps College,<sup>2</sup>Bucknell University, <sup>3</sup>Lawrence Berkeley National Laboratory

### **Research Question**

How well can multidimensional similarity based compression work on scientific data?

# **Dynamic Time Warp**

•DTW performs nonlinear "warping" on the sequences where differences in time are not penalized

•For time series of length n, it is necessary to do n<sup>2</sup> computations

 $d_{\mathrm{DTW}}(\mathbf{x},\mathbf{y}) = D_{M,N}$  $D_{i,j} = d_{\text{Euc}}(x_i, y_i) + \min\{D_{i,j-1}, D_{i-1,j}, D_{i-1,j-1}\}$ 

$\begin{bmatrix} 3 \\ 3 \\ 4 \\ 5 \\ 3 \end{bmatrix}$	1 1 2 3 1	$2 \\ 2 \\ 3 \\ 4 \\ 2$	$egin{array}{c} 0 \\ 0 \\ 1 \\ 2 \\ 0 \end{array}$	$1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 1$	$egin{array}{c} 0 \\ 0 \\ 1 \\ 2 \\ 0 \end{array}$	$D(\mathbf{x}, \mathbf{y}) =$	$ \begin{bmatrix} 3 \\ 6 \\ 10 \\ 15 \\ 18 \end{bmatrix} $	- 4 - 4 - 4 - 6 - 9 - 10	$     \begin{array}{r}       - 6 \\       6 \\       7 \\       10 \\       11     \end{array} $	$     \begin{array}{c}       6 \\       6 \\       7 \\       12 \\       10     \end{array} $	7 7 6 7 8	7 7 7 8 7	
$\begin{vmatrix} 0\\3\\7 \end{vmatrix}$	$\frac{1}{5}$	2 6		$\frac{1}{3}$	$\begin{vmatrix} 2\\0\\4 \end{vmatrix}$		18 25	10 15	11 16	10 14	8 11	7 11	>

 $d_{\mathrm{DTW}} = D_{6,6}$ 

= 11

### Minimum Jump Cost

•MJC works by accumulating the cost of jumping forward from one time series data point to the nearest data point in the other time series  $d_{MJC} = \sum c_{\min}^{(i)}$  $c_{\min}^{(i)} = \min\{c_{t_x}^{t_y i}, c_{t_x}^{t_y+1}, c_{t_x}^{t_y+2}, \ldots\}$ 

Figure 4: MJC for the first data point (left). Total jumps (right).

• Instead of calculating all n<sup>2</sup> distance values of between **x** and **y**, only the distance between points of index greater than the recursive starting point are calculated

• Expected to reduce runtime







MJC has lower error at larger dictionary size

Table 1: Dictionary size comparison for 100 CR



**SCRIPPS** 

THE WOMEN'S COLLE

computationally expensive