

# 基于序列相似性计算的甲骨残片缀合算法

张重生<sup>1,2</sup>, 王斌<sup>1</sup>

(1. 河南大学河南省大数据分析与管理重点实验室, 河南开封 475001; 2. 河南大学黄河文化遗产实验室, 河南开封 475001)

**摘要:** 甲骨残片缀合一直是甲骨学研究中最急迫最具基础性的工作, 它使得甲骨残片经过拼接, 复原为更加完整的原始材料. 尽管前人及同行曾提出若干计算机辅助的甲骨缀合方法, 但这些方法缀合准确度不足, 未能真正投入使用, 并不能真正帮助专家解决甲骨缀合问题, 导致当前的甲骨缀合工作仍旧依靠人工、依旧费时费力. 为了更好地研究甲骨残片的机器缀合问题, 本文使用一个较大规模甲骨缀合基准数据集 OB-Rejoin, 该数据集包含了约一千幅甲骨拓片图像, 且融入了大量的甲骨学界已缀成果, 用于算法评估. 基于该数据集, 本文设计了一种基于斜率变化量序列匹配的甲骨缀合算法 (Slope United Sequence Matching for Oracle Bone Fragments Conjugation, SUM), 该方法将甲骨残片的断边磕口图像匹配问题转化为数值型的序列数据和序列相似性比对问题, 以将尚不够非常精密的计算机视觉领域的磕口图像匹配问题转换为数据科学领域较为成熟的序列数据相似性匹配问题. SUM 将数值型的磕口序列数据进一步转换为斜率变化量序列和字符序列数据, 最后利用字符序列的模糊匹配完成甲骨残片的磕口匹配. 在实验环节, SUM 算法与经典的序列相似性计算方法在精确率、召回率、漏检率方面进行了对比, 并与两个较新的基于深度学习的序列匹配算法和形状匹配算法进行了性能对比. 整体而言, SUM 在 OB-Rejoin 数据集上的 Top-15 缀合召回率达到了 95.181%, 超越了对比算法. 重要出土文献的精准复原本身是历史学和古文字研究中客观存在的重大现实需求, 具有重要的史学价值和意义, 因此, 本文的研究成果, 不但有助于解决甲骨残片的机器缀合问题, 还对秦汉简牍和敦煌遗书等重要出土文献的精准复原具有重要的参考价值.

**关键词:** 甲骨文; 甲骨缀合; 序列相似性计算; 形状匹配; 边缘匹配

**基金项目:** 科技部高端外国专家项目 (No.G2021026016L)

**中图分类号:** TP181

**文献标识码:** A

**文章编号:** 0372-2112(2023)04-0860-10

**电子学报 URL:** <http://www.ejournal.org.cn>

**DOI:** 10.12263/DZXB.20210429

## Oracle Bone Fragments Conjugation Based on Sequence Matching

ZHANG Chong-sheng<sup>1,2</sup>, WANG Bin<sup>1</sup>

(1. Henan Key Laboratory of Big Data Analysis and Processing, Henan University, Kaifeng, Henan 475001, China;

2. Laboratory of the Yellow River Heritage, Henan University, Kaifeng, Henan 475001, China)

**Abstract:** Rejoining the oracle bone fragments is an important prerequisite for the research of oracle bone inscriptions (OBI), which can restore the original appearance and content of the oracle bones. Though computer-aided oracle bone fragments conjugation solutions have been investigated for decades, they could not be applied in real-world OBI research, due to their unsatisfactory performance. Consequently, until today, OBI researchers still have to rejoin the oracle bone fragments manually. To solve this problem, we first introduce OB-Rejoin, a large-scale dataset with about one thousand oracle bone rubbings. It includes a large number of fragments that have already been rejoined by OBI experts, which are used as the ground-truth in experiments. Moreover, we propose the SUM (Slope United Sequence Matching) algorithm for oracle bone fragments conjugation, which transforms the challenging curve matching problem of the oracle bone fragments into the numerical sequence matching problem. SUM next transforms the sequence data into slope variation-based sequence data and character sequences, and finally uses string matching algorithms for oracle bone fragments conjugation. We conduct comprehensive experiments to compare SUM with classic sequence matching methods, in terms of precision, recall, mis-rejoin rates. We also compare SUM with two very recent deep learning-based sequence matching and shape matching algorithms. All these experiments demonstrate the superiority of SUM over existing methods in oracle bone fragments conjugation, which achieves a Top-15 recall rate of 95.181% on OB-Rejoin. Overall, the recovery of unearthed docu-

ments is an important real-world problem that has historical significance, this research work is therefore not only useful for rejoining the oracle bone fragments, but also has important reference value for the recovery of other unearthed documents, in particular the conjugation of fragmented bamboo strips and Dunhuang manuscripts.

**Key words:** oracle bone inscriptions; oracle bone fragments conjugation; sequence matching; shape matching; edge matching

**Foundation Item(s):** High-end Foreign Expert Project of the Ministry of Science and Technology of China (No. G2021026016L)

## 1 引言

甲骨文是我国目前可见的最早成体系的文字,是中华民族珍贵的文化遗产。龟甲和兽骨等材料是甲骨文的主要载体,而甲骨出土和甲骨缀合是获得甲骨学新材料的两种主要方法<sup>[1]</sup>。甲骨之所以需要缀合,是因为出于占卜的需要,古人对龟甲和兽骨进行了钻凿和灼烧,使其变得容易断裂,又因埋藏地下 3 000 余年,受到长期腐蚀而残碎;又因早期考古发掘手段及战争年代的运输和流传等因素的影响,变得更加残损。只有将这些残碎的甲骨缀合起来,才能提供更为全面的卜辞内容<sup>[2]</sup>。因此,甲骨缀合一直是甲骨学研究中最具基础性的工作,得到了甲骨学界的持续投入和关注。

然而,目前的甲骨缀合主要依靠人工,如刘影<sup>[3]</sup>、齐航福<sup>[4]</sup>根据上下文语义和甲骨形态学等方面的专业知识进行手工缀合,虽取得了较好的缀合成果,但时间消耗巨大、工作效率较低。随着计算机技术的发展,也有一些计算机辅助的甲骨缀合尝试,但停留在研究的初级阶段。王爱民等人<sup>[2,5]</sup>尝试利用边缘匹配方法进行甲骨缀合,但所用的实验数据较少,准确度不够理想,没有真正运用在实际的甲骨缀合研究中。

为了更好地开展甲骨残片的机器缀合研究,本文使用了一个新的甲骨缀合数据集 OB-Rejoin。该数据集利用甲骨断边碴口匹配(边缘匹配)的方法进行甲骨缀合,同时融入了缀合后图像的边缘平滑性等甲骨形态学知识,减少缀合干扰项,提高缀合精度,减轻专家筛选缀合结果的工作负荷。具体而言,甲骨专家使用红色曲线描绘断裂边缘,使用绿色直线表示甲骨原始边缘(非断裂形成的边缘,简称原边)的弧度走向。根据甲骨形态学的知识判断缺失部位,将 OB-Rejoin 数据集分为两部分,即上部拓片和下部拓片,各包括 499 幅图像,共计 998 幅图像。该数据集还融入了甲骨专家已成功缀合的 249 对图像,以便为缀合算法提供参照,检验其有效性。

基于 OB-Rejoin 数据集,本文设计 SUM 甲骨缀合算法(Slope United Sequence Matching for Oracle Bone Fragments Conjugation),即基于斜率变化量序列匹配的甲骨缀合算法,简称 SUM 算法。该算法首先通过可回溯搜索算法提取甲骨断边碴口图像对应的数值型坐标序列

数据,实现甲骨断边碴口图像到坐标序列数据的转换。接着,将该坐标序列数据进一步转换为斜率变化量序列和字符序列,最终通过字符序列之间的相似性匹配实现甲骨断边碴口图像之间的匹配,完成甲骨残片的机器缀合。

为了验证 SUM 算法的有效性,本文在基准数据集 OB-Rejoin 上,将 SUM 算法与代表性的序列相似性计算方法在精确率、召回率、漏检率方面进行大规模实验对比分析,其 Top-15 缀合召回率达到了 95.181%,验证了 SUM 甲骨缀合算法的优越性。本文还将 SUM 算法与近年基于深度学习的序列匹配算法和形状匹配算法进行了对比分析,进一步说明了其优势。

本文的主要贡献概括如下:

(1)使用了一个较大规模的 OB-Rejoin 数据集作为甲骨缀合的基准数据集,该数据集共计 998 幅甲骨拓本图像且融入了甲骨学领域知识,使用红色曲线标注甲骨断边碴口,使用绿色直线表征甲骨原边的弧度走向特征。

(2)提出将断边碴口图像匹配问题转换为数值型序列数据匹配问题的甲骨缀合新思路,设计了一种新颖的基于序列计算的缀合方法 SUM,将甲骨断边碴口边缘图像转换为数值型的坐标序列数据,再进一步转换为斜率变化量序列和字符序列,最终通过字符序列模式匹配实现甲骨残片的机器缀合。

(3)基于 OB-Rejoin 数据集上开展了大规模实验,证明了 SUM 算法较 DTW 算法在召回率方面有显著提升;而且,与近年的基于深度学习序列匹配算法和形状匹配算法相比,SUM 仍有明显优势。

## 2 相关工作

甲骨拓片图像是甲骨学研究的主要材料。甲骨缀合可分为两种方式,一种是纯手工的专家缀合甲骨方式<sup>[3,4,6]</sup>,另一种是有计算机辅助、专家参与的甲骨缀合方式<sup>[5-11]</sup>。

专家缀合甲骨的方式,通常由甲骨专家综合利用上下文语义、可以拼接的文字、文献资料及甲骨学形态等知识推测候选缀合结果。但由于人力的局限性,缀合效率较低。尤其在面对大量的甲骨碎片时,专家缀合甲骨的方式会陷入困境。甲骨学界一直对利用计算机技

术实现甲骨缀合寄予厚望,已有计算机辅助缀合方法主要是基于甲骨碴口图像边缘匹配的方法,首先检测碴口边缘,再通过边缘特征比对,选择候选缀合结果.但这些方法的缀合准确度较低<sup>[5-11]</sup>,未投入实际使用.

时间序列分析理论<sup>[12]</sup>已引起学术界和工业界广泛关注,在医药、金融、科学研究领域具有广泛应用,而DTW(Dynamic Time Warping)算法是该领域最经典、最有效的方法之一.下面将DTW算法的相关研究成果进行系统归纳,如表1所示.

表1 DTW相关的论文概要

参考文献	年份	主要内容
文献[13]	2016	纠正了某些文献中“DTW效果不佳”或者“DTW的窗口大小无关紧要”等误解,总结DTW算法性能优化方面的经验.
文献[14]	2000	作者证明常用的欧氏距离,是一个非常脆弱的距离度量方法,存在时间轴畸变时退化迅速的情况.作者演示了DTW的一种改进方法,使用时间序列数据的更高级别表示,大大提升了计算效率.
文献[15]	2011	基于两个序列之间的形状相似性是准确识别的主要因素,提出一种新的距离测度加权DTW(WDTW),用来解决DTW不能说明参考点和测试点之间相位差的相对重要性可能导致的错误分类.该方法对参考点和测试点的相位差较大的点进行惩罚,以防止异常值引起的最小距离畸变.实验结果表明,该方法可以提高时间序列分类分析的精度.
文献[16]	2005	指出基于DTW的时间序列比对研究中的一些经验总结: (1) 比较不同长度的序列并将它们重新插补到相同长度,DTW在精度/召回率方面不会产生显著差异. (2) 为了让DTW易于处理,约束扭曲路径并不是一个好办法. (3) 试图通过产生更严格的下限来加速DTW的计算时间几乎是毫无意义的.
文献[17]	2017	DTW得到的是全局最优解,而不是局部敏感匹配.为了让DTW匹配局部结构完全不同的两个时间点,作者提出了一种改进的对齐算法——shapeDTW,通过考虑局部结构信息,提高了DTW的性能.
文献[18]	2019	作者提出了一种基于人工神经网络的DTW实现方式,以获得更好的特征提取和匹配准确度.
文献[19]	2016	提出一个精确加速所有成对DTW矩阵计算的方法,平均减少了大约50%的运行时间.
文献[20]	2017	在DTW中对翘曲路径施加约束,避免可能导致的错误对齐.
文献[21]	2019	证明DTW相似度不是扭曲不变的,并将其转换为扭曲不变的度量方法,该方法需要更少的存储和计算时间.
文献[22]	2007	验证DTW在领域无关和领域相关的时间序列分类方面的性能差异.发现在训练数据较少时,使用领域相关的算法比较合适;在一些形状比较规则的数据集上,使用DTW方法有一定的优势.
文献[23]	2020	综述了近年时间序列分析在数据挖掘领域的文献成果,对时间序列特征表示和相似性度量方法进行了总结.
文献[24]	2018	基于DTW方法进行时序数据相似连接,即寻找给定相似度度量下的所有相似时序数据对,设计了两种高效计算方法.
文献[25]	2009	基于DTW和HMM,进行语音识别.
文献[26]	2012	作者介绍了基于轮廓和骨架的形状表示与匹配的分类方法,同时介绍和分析了形状检索中的度量学习.
文献[27]	2007	作者提出了一种线性时间复杂度和空间复杂度的改进DTW算法.
文献[28]	2022	作者介绍了有关计算时间序列相似度时的一些优化技巧.

### 3 算法设计

#### 3.1 问题定义

甲骨缀合是甲骨学研究中客观存在的重要现实需求.通过计算机技术实现甲骨残片的自动缀合,具有重要的史学价值和意义.本文拟解决的问题为:对每幅甲骨图像,计算与其断边碴口图像最为匹配的前 $K$ 幅图像,并保证其中包含能真正缀合的结果.

#### 3.2 算法整体框架

所设计的SUM算法(Slope United sequence Matching for oracle bone fragments conjugation)的整体框架如图1所示,主要包含三部分/三个步骤:(1)首先将每幅甲骨残片的断边碴口图像转换为数值型的坐标序列数据;(2)基于该坐标序列数据,再进一步转换为斜率变

化量序列和字符序列数据;(3)设计/使用字符序列相似性计算方法,通过字符序列的模糊匹配完成甲骨残片的机器缀合.

#### 3.3 提取甲骨断边碴口图像对应的坐标序列数据

对于每幅拓片图像,设计算法提取其断边坐标序列.使用深度优先搜索的方法,搜索条件设置为一个像素点相邻的八个方向,需满足八个方向至少有一个方向出现红色像素点.然后,再由搜索到的像素点出发,继续搜索下一个红色像素点的坐标,直到围绕断边一周,形成环路,搜索结束.

#### 3.4 计算斜率梯度序列和斜率变化量序列

对于每幅拓片图像,设计算法提取其断边坐标序列.使用深度优先搜索的方法,搜索条件设置为一个像

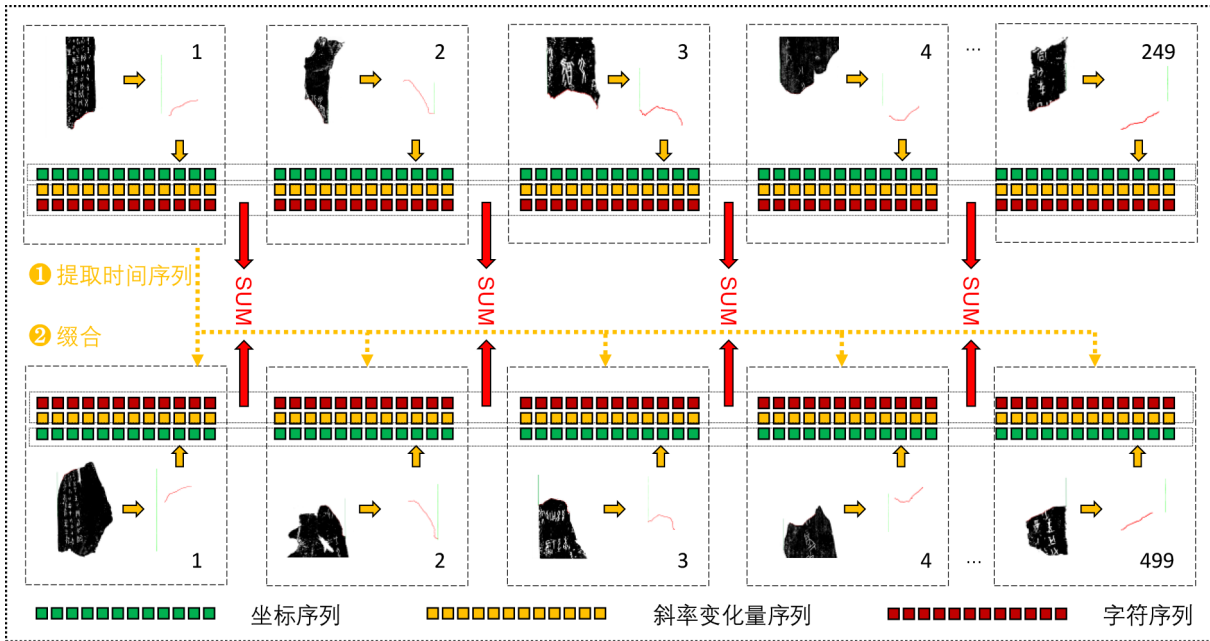


图1 SUM甲骨缀合算法的整体计算流程

素点相邻的八个方向,需满足八个方向至少有一个方向出现红色像素点. 然后,再由搜索到的像素点出发,继续搜索下一个红色像素点的坐标,直到围绕断边一周,形成环路,搜索结束.

得到甲骨断边碴口图像对应的数值型坐标序列后,SUM甲骨缀合算法设计邻位斜率法和尾部斜率法两种方法,计算斜率梯度序列. 邻位斜率法对所有两两相邻的点直接求斜率,但该方法只对波形敏感,而对振幅不敏感. 为了解决该问题,我们设计了尾部斜率法,其计算方法如下,令坐标序列  $A = \{(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)\}$ , 其中  $n$  是序列  $A$  的长度,  $(x_i, y_i)$  是某个像素点的坐标值. 计算斜率梯度序列  $A' = \{f'_{x_1}, f'_{x_2} \dots f'_{x_n}\}$ , 公式如式(1)所示:

$$f'_{x_i} = \frac{y_{i'} - y_i}{x_{i'} - x_i} (y_{i'} \neq y_i, x_{i'} \neq x_i, i + 5 \leq n) \quad (1)$$

其中,  $f'_{x_i}$  是  $x_i$  处的斜率,我们将每一个点与它后临的若干个点(本文根据经验值,将该参数设置为5个点)的坐标依次比对,如果五次均出现两个点同行或同列,将  $f'_{x_i}$  置为1;否则,取点  $i$  后邻的五个点中最后一个符合与点  $i$  既不同行又不同列的点,如式(1)所示计算点  $i$  的斜率  $f'_{x_i}$ ,即纵坐标的差除以横坐标的差.

在求得斜率梯度序列的基础上,相邻的斜率值逐位相减,得到斜率变化量序列. 即,某个点的斜率变化量为下一点的斜率值减去该点的斜率值. 然后,根据每个斜率变化量的值所在的值域区间,将斜率变化量序列进一步转换为字符序列.

图2是上述计算过程的示例:SUM甲骨缀合算法对

提取到的断边坐标序列,首先使用尾部斜率法计算斜率梯度序列,再由斜率梯度序列得到斜率变化量序列和其抽象后的字符序列.

### 3.5 字符序列匹配算法

甲骨残片的断边普遍存在部分残损、边缘粗糙、特征模糊等情况,因此,很少有断边完全密合的甲骨缀合结果. 所以,缀合算法需要支持甲骨断边间的模糊匹配,SUM甲骨缀合算法通过字符序列的模糊匹配完成甲骨缀合. 在进行字符序列的模糊匹配时,考虑了上部序列优先和下部序列优先两种缀合方式.

图3以示例形式阐述了字符序列匹配算法的计算过程. 具体计算过程如下:

Step 1 初始时,选取上部字符序列第1个字符和下部字符序列第1个字符.

Step 2 对于上部字符序列的当前字符,与下部字符序列的字符开始比对,移动下部序列的下标,在比对不超过6次找到相同字符,进行Step 4;否则,下标复位,进行Step 3.

Step 3 选中下部字符序列的当前字符,与上部字符序列的字符开始比对,移动上部序列的下标,在比对不超过6次找到相同字符进行Step 4,否则将上部和下部字符序列的下标同时后移1位,继续Step 2.

Step 4 上下部字符序列从第1对可以匹配的字符开始,依次向后比对,同时移动上下部的下标,若相同的字符不超过3个,两者下标复位后同时后移一位,返回Step 2. 否则,记录相同的个数,更新两者下标,返回Step 2.

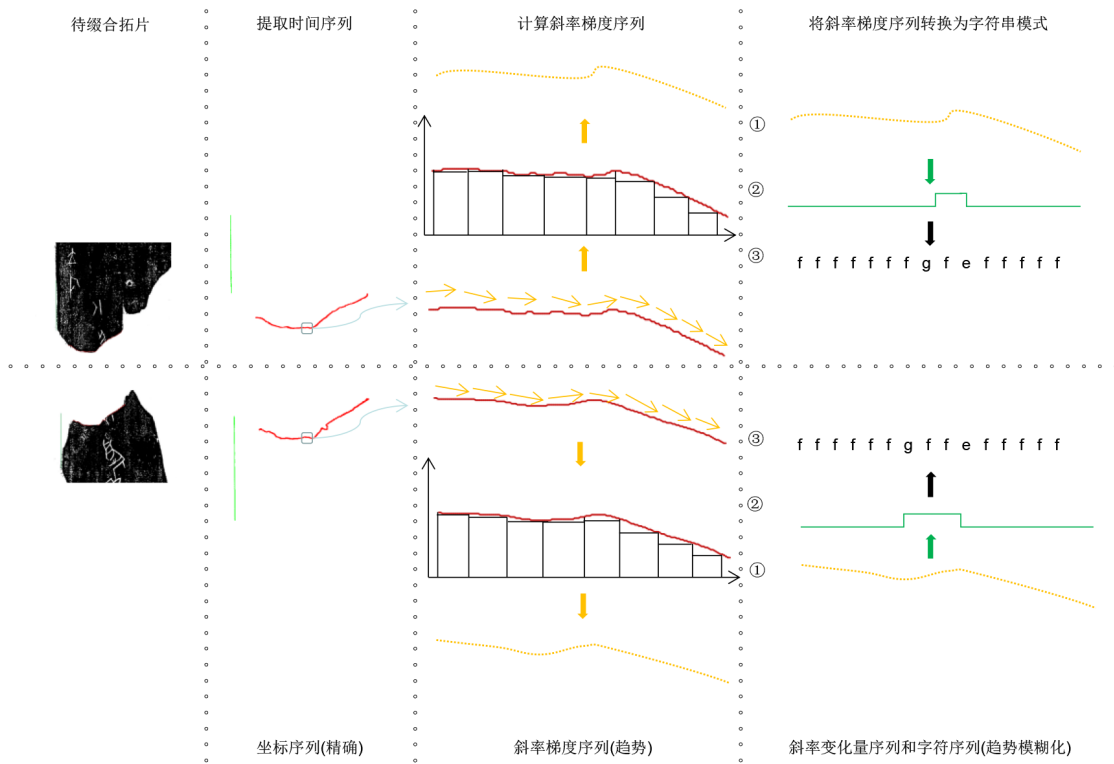


图2 斜率变化量序列和字符序列的计算过程

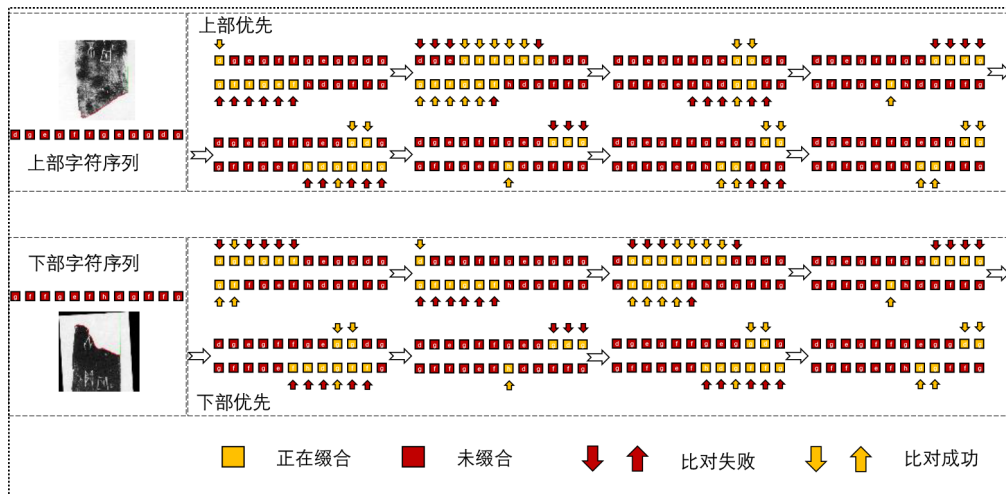


图3 SUM算法的字符匹配步骤示例

上述比对结束的条件是任何一个序列的下标搜索到达终点。

### 3.6 加入甲骨形态学领域知识的必要性

甲骨专家使用红色曲线,描绘断边的边缘;使用绿色直线,指示甲骨原边(非断裂形成的边缘)的弧度走向.甲骨原边的弧度走向属于甲骨形态学知识,来自甲骨学家的建议.如果只有断边匹配、没有原边弧度走向一致性的判断,则很难保证缀合后的原边是否光滑.而甲骨学家在人工缀合时,通常利用原边弧度走向是否一致性剔除干扰性的缀合结果.因此,本文在首先保证

两幅待缀甲骨的原边弧度走向一致的前提下,再计算断边的边缘相似度,以降低计算量,减少干扰性的候选缀合结果.本文在开展相关实验时,为保证公平性,统一对相关对比算法使用了原边弧度走向的预处理步骤.

## 4 实验结果与分析

### 4.1 实验设置

**实验数据集** 本文实验所用的基准数据集为OB-Rejoin数据集.该数据集包括499幅上部拓片和499幅下部拓片,并在其中融入了249对甲骨专家通过

人工成功缀合的甲骨拓片图像(分别放置在上部拓片和下部拓片的文件夹中)。

上部拓片和下部拓片是指甲骨专家根据缺失部位判断每幅图像为上部拓片或下部拓片。若为上部拓片,则该图像的断边在下部,即缺失部分在下部;若为下部拓片,则该图像的断边在图像上部,即缺失部分在上部。两幅图像若能成功缀合,则缀合后上部图像在上、下部图像在下。简言之,上部拓片文件夹共有 499 幅甲骨图像,其中的 249 幅在下部拓片文件夹中存在能与之真正缀合的甲骨图像。

对每幅甲骨图像,甲骨专家手工描绘甲骨断边的边缘图像(红线)及表示甲骨原边(侧边)平滑弧度走向的直线(绿线),其表现形式如图 4 所示。甲骨专家设计绿线特征的原因是:实际的甲骨缀合工作对缀合算法的准确度有极高的要求,例如,在实验数据集上 Top-K 召回率(下面将解释该指标的含义)为 80% 左右的缀合算法,根本无法在实际缀合工作中应用,因为实际的缀合工作难度更大、干扰项更多。而甲骨专家手工缀合甲骨时,他们通常会查看两幅甲骨拓片图像试缀后两幅图像的原边(侧边)的弧度走向是否一致,用以排除干扰性的候选结果。因此,甲骨专家建议设计表示甲骨原边平滑弧度走向的特征(绿线),以帮助算法进一步提升缀合效果。相关的消融实验在 4.2.2 节给出。

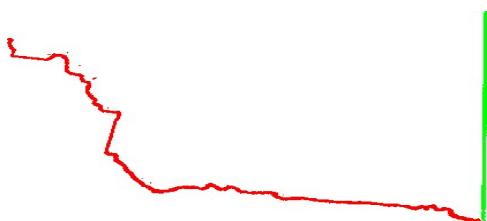


图 4 从 OB-Rejoin 数据集中的每幅甲骨拓片图像提取到的断边边缘(红线)和原边弧度走向特征(绿线)效果图

**评估指标** 考虑到甲骨缀合的实际应用场景为:输入一幅甲骨残片图像,返回与该图像断边磕口(边缘)匹配度最高的前  $K$  (Top- $K$ ) 幅甲骨图像作为候选缀合结果,并期望  $K$  幅图像中包含能与输入图像真正缀合的甲骨残片图像,因此,本文的主要评估指标是 Top- $K$  召回率/查全率,即为每幅图像返回缀合匹配度最高的前  $K$  幅候选图像的情况下,能正确找到缀合结果的图像数量与实际有缀合对象的全部图像数量之比。需要说明的是,有些甲骨残片图像可能不存在真正能缀合的对象。如在 OB-Rejoin 基准数据集中,上部图像共计 499 幅,其中只有 249 幅图像存在能真正缀合的对象,而其余 250 幅图像并不存在真正能缀合的对象。在实验数据集 OB-Rejoin 上,每幅图像的真实缀合对象已知,缀合算法可以自动计算其 Top- $K$  召回率,无须专家参与。但在真实的甲骨缀合研究的场景中,由于事先并不知

道某幅图像的真实缀合对象,故需要甲骨专家通过专业知识进行筛选,确认候选缀合结果。若  $K$  值很大,则候选结果过多,查看每幅甲骨残片图像的候选结果将耗费甲骨专家大量的时间;又由于甲骨图像数量庞大,甲骨学者将很难完成所有甲骨残片图像候选缀合结果的筛选和确认。为了切实减轻真实甲骨缀合场景下甲骨学者的工作负荷,结合甲骨专家的建议,将  $K$  值设置为 1、3、5、10、15 等取值。

除了 Top- $K$  召回率外,本文还考虑了漏检率、精度等指标。其中,漏检率=1-Top- $K$  召回率;而精度有两种计算方式:精度 1 指正确缀合的图像数量除以候选结果总量(每幅图像均有  $K$  个候选结果),精度 2 指正确缀合数量除以对应的真正缀合结果的在  $K$  个候选结果中的排名之和。若两个缀合算法寻找到的正确缀合结果的数量相同,则其精度 1 相同,Top- $K$  召回率和漏检率亦相同;此种情况下,若其中的某个方法寻找到的正确缀合结果在  $K$  个结果中的排名均靠前,则其精度 2 值更大。若某个方法寻找到的正确缀合结果在  $K$  个结果中均排名第一 (Top-1),则其精度 2 指标值为 100%。总体而言,高的 Top- $K$  召回率和高的精度指标值(含精度 1、精度 2)是甲骨缀合算法的期望效果。

**对比算法** 本文工作的相关方法主要涉及序列匹配、形状匹配/轮廓匹配等方面的算法<sup>[29]</sup>。因此,除了 DTW<sup>[15]</sup> 和 Discrete Fréchet distance<sup>[29]</sup> 等经典的序列匹配算法之外,本文还增加了基于深度学习的序列匹配算法 DTWNet<sup>[18]</sup>,代表性的形状匹配算法 Shape Context Matching 算法<sup>[30]</sup> 和较新的基于深度学习的形状匹配算法 Deep Shape Matching<sup>[31]</sup> 等对比算法。按方法类型,DTW、DTWNet 和 Discrete Fréchet distance 属于序列匹配算法,而 Shape Context Matching 和 Deep Shape Matching 属于形状匹配/轮廓匹配算法。尤其值得一提的是,DTWNet 和 Deep Shape Matching 分别为近年提出的基于深度学习的序列匹配算法和形状匹配算法。有关 DTW 算法,本文使用文献[15]对应的 DTW 算法实现,通过在 OB-Rejoin 数据集上的实验尝试,发现 DTW 算法的窗口值设置为 85 时效果最好;我们还在 DTW 中加入了归一化处理,以减小误差,提升其在 OB-Rejoin 数据集上的算法性能。

## 4.2 实验结果分析

### 4.2.1 对比实验

#### (1) Top- $K$ 召回率评估指标

如前文所述,OB-Rejoin 数据集包括上部拓片文件夹和下部拓片文件夹两部分,各含 499 幅图像;上部拓片文件夹中的 249 幅甲骨图像,在下部文件夹中存在能与之真正缀合的对象。对上部拓片文件夹中的每幅图像,算法须在下部拓片文件夹中寻找与其缀合匹配度

最高的前  $K$  幅图像. 评估指标主要使用 Top- $K$  召回率/查全率,  $K$  值取 1、3、5、10、15.

用于甲骨缀合对比实验的算法包括 DTW<sup>[15]</sup>、DTWNet<sup>[18]</sup>, Discrete Fréchet distance<sup>[29]</sup>, Shape Context Matching<sup>[30]</sup>和 Deep Shape Matching<sup>[31]</sup>五个方法, 其中的前三个方法为序列匹配算法, 后两个方法为形状匹配算法. 本实验使用的对比算法源代码均为公开的算法实现, 通过大量的程序调试和适配, 将其应用于甲骨缀合问题上.

表 2 给出了本文所设计的 SUM 算法及五个对比算法在 OB-Rejoin 数据集上的缀合性能, 基于 Top- $K$  召回率/查全率评估指标. 从表中可以看出, 在三个序列匹配算法中, Discrete Fréchet distance<sup>[29]</sup>方法的缀合结果显著优于其它两个方法, 含基于深度学习的 DTWNet 方法(已使用了大量的数据增强技术); 而在两个形状匹配算法中, 基于深度学习的 Deep Shape Matching<sup>[31]</sup>方法的性能较好, 且超越了三个序列匹配算法. 相较于 Deep Shape Matching 方法, SUM 算法在 Top-1 召回率指标上的提升接近 10%, 非常显著; 在 Top-3 召回率指标、Top-5 召回率指标上、Top-10 召回率指标上、Top-15 召回

率指标上, SUM 算法的性能提升在 2.4%~4% 之间, 提升明显. 整体而言, 本文所设计的 SUM 算法取得了最佳的缀合效果, 尤其是 Top-15 召回率达到了 95.181%, 体现了该算法的优越性.

表 2 SUM 方法与对比算法的 Top- $K$  召回率/查全率性能对比, 含 DTW<sup>[15]</sup>, DTWNet<sup>[18]</sup>, Discrete Fréchet distance<sup>[29]</sup>, Shape Context Matching<sup>[30]</sup>, Deep Shape Matching<sup>[31]</sup>等 5 个方法的缀合准确率 单位:%

Top- $K$ 召回率	文献 [15]	文献 [18]	文献 [29]	文献 [30]	文献 [31]	SUM
Top-1	33.330	38.554	57.831	47.791	60.643	70.281
Top-3	46.988	55.422	76.707	67.068	79.518	82.731
Top-5	55.820	65.462	82.329	75.100	85.542	88.353
Top-10	65.860	74.297	87.149	82.329	88.755	92.771
Top-15	72.290	83.133	91.165	86.345	92.771	95.181

## (2) 其他评估指标

表 3 展示了 DTW 算法和 SUM 算法在精度、召回率、漏检率方面的性能对比, 其中的 DTW 指没有利用甲骨原边弧度走向一致性特征(简称绿线特征), 而预处理下的 DTW 指利用绿线特征的 DTW 算法. 整体而言, SUM 算法在各个指标上的性能均显著优于 DTW 算法.

表 3 DTW、预处理下 DTW 和 SUM 算法的缀合准确率

单位:%

算法	DTW	预处理下的 DTW	SUM
精度 1	Top-1: 20.48	Top-1: 33.33	Top-1: 70.28
	Top-5: 7.39	Top-5: 11.16	Top-5: 17.67
	Top-10: 4.50	Top-10: 6.59	Top-10: 9.28
	Top-15: 3.27	Top-15: 4.820	Top-15: 6.35
精度 2	Top-1: 20.48	Top-1: 33.33	Top-1: 70.28
	Top-5: 9.52	Top-5: 17.10	Top-5: 48.46
	Top-10: 6.56	Top-10: 12.59	Top-10: 40.38
	Top-15: 5.16	Top-15: 10.61	Top-15: 37.19
召回率	Top-1: 20.48	Top-1: 33.33	Top-1: 70.28
	Top-5: 36.95	Top-5: 55.82	Top-5: 88.35
	Top-10: 44.98	Top-10: 65.86	Top-10: 92.77
	Top-15: 49.00	Top-15: 72.29	Top-15: 95.18
漏检率	Top-1: 79.52	Top-1: 66.67	Top-1: 29.72
	Top-5: 63.05	Top-5: 44.18	Top-5: 21.65
	Top-10: 55.02	Top-10: 34.14	Top-10: 7.23
	Top-15: 51.00	Top-15: 27.71	Top-15: 4.82

## 4.2.2 消融实验

### (1) 加入甲骨形态学知识对缀合结果的影响

在实验过程中, 基于真正缀合的两幅甲骨图像的原边弧度走向一致的甲骨形态学领域知识, 首先将试缀的两幅甲骨图像均按绿线方向竖直放置, 以保证两幅甲骨缀合后的原边弧度走向一致. 该特征有利于排除一部分干扰性缀合结果. 需要说明的是, 在进行上述综合实验对比分析时, 所有方法均使用了绿线特征, 即

所有方法都对每幅甲骨图像沿绿线方向进行了竖直放置. 为了说明该特征对甲骨缀合的作用, 本部分消融实验中, 尝试去掉绿线特征, 对上部图像和下部图像各进行一定角度的旋转( $3^\circ$ ,  $6^\circ$ ), 对应的 SUM 算法的实验结果如图 5 所示. 可以观察到, 使用绿线特征对 SUM 算法的缀合效果有一定的提升作用, 说明了加入甲骨形态学领域的特征有助于甲骨缀合效果的提升.

另外, 在 4.2.1 节实验中, 表 3 对比了是否利用甲骨

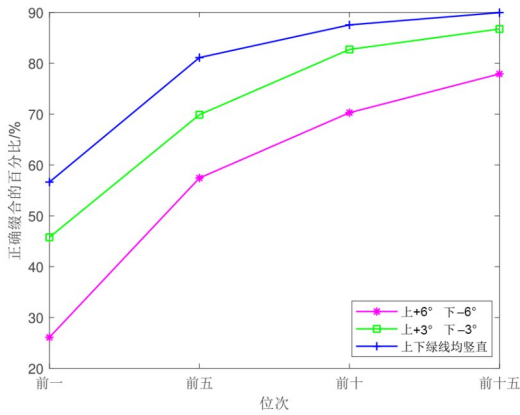


图5 加入甲骨形态学领域知识对缀合结果的影响

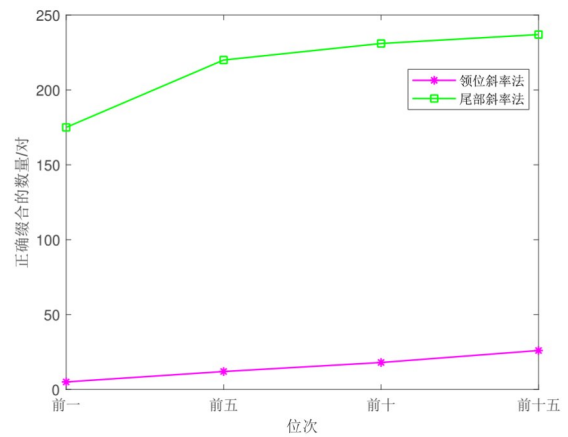


图6 邻位斜率法和尾部斜率法缀合效果对比

原边弧度走向一致性特征的两种 DTW 方法的性能, 对应结果进一步说明了甲骨原边弧度走向一致性特征 (绿线特征) 对甲骨缀合算法的帮助.

(2) 两种斜率梯度序列计算方法的性能对比

由于甲骨断边磕口 (边缘) 并不平滑, 因此对应的转换后的数值型序列数据包含了大量干扰信息. 而邻位斜率法极易受到干扰信息的影响, 为消除该影响, 本文设计了对波形和振幅均敏感的尾部斜率法. 所使用的两种斜率梯度序列计算方法对应的缀合结果如图6所示. 从图中可以看出, 尾部斜率法具有更好的缀合效果. 因此, 在相关实验中, SUM 算法均采用尾部斜率法计算斜率梯度序列.

4.2.3 甲骨缀合结果展示

在 SUM 算法的缀合结果中, 有 175 幅甲骨图像的 Top-1 缀合结果即为真正缀合结果. 图7展示了部分 Top-1 缀合结果示例. 其 Top-10 缀合结果如图8所示, 对应的甲骨磕口边缘曲线的整体轮廓更为平直, 特征越发不明显; 磕口边缘曲线也出现了较多反弓, 缀合难度逐渐增大.

图9展示了一个失败的甲骨缀合案例. 由于输入的甲骨图像磕口边缘短直、特征点较少, 导致能匹配的对象较多, 使得 Top-K 缀合结果不能包含能与之真正缀合的甲骨图像, 导致缀合失败.

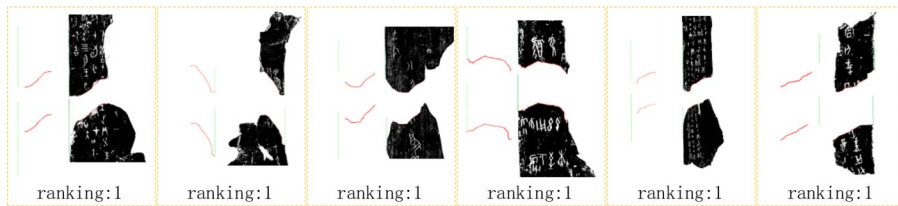


图7 SUM算法的 Top-1 缀合结果示例

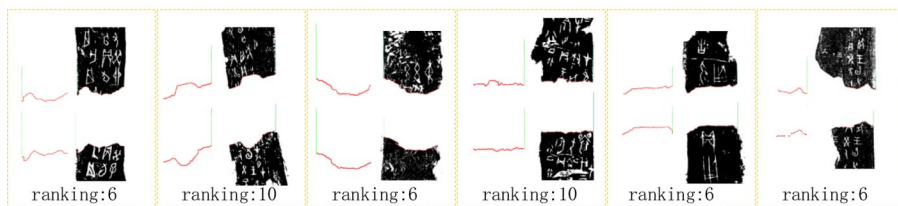


图8 SUM算法的 Top-10 缀合结果示例



图9 缀合失败案例

## 5 结束语

为研究甲骨残片的机器缀合问题,本文使用了一个较大规模的基准数据集 OB-Rejoin,设计了甲骨缀合算法 SUM,该算法将甲骨碴口边缘图像转换为坐标序列数据,并进一步转换为斜率变化量数据和字符序列数据,最后通过字符序列的模糊匹配实现甲骨残片的机器缀合.在基准数据集上,SUM算法的 Top-1 缀合召回率为 70.281%,Top-5、Top-10 和 Top-15 缀合召回率分别达到了 88.353%、92.771% 和 95.181%,缀合效果较为显著.总之,SUM 算法的主要创新在于将计算机视觉领域的甲骨图像边缘匹配问题转化为数据科学领域的序列数据匹配问题的新思路,而非具体的字符序列匹配算法设计.本文的研究成果,对我国重要出土文献的精准复原具有一定的参考价值,对古文字研究具有一定的推动作用.

## 参考文献

- [1] 黄天树. 甲骨缀合的学术意义与方法[J]. 故宫博物院院刊, 2011, 54(1): 7-13, 156.  
HUANG T S. On academic value and research methods of restoration of oracle bone inscription fragments[J]. Palace Museum Journal, 2011, 54(1): 7-13, 156. (in Chinese)
- [2] 王爱民, 葛文英, 赵哲, 等. 龟甲类甲骨文碎片计算机辅助缀合研究[J]. 计算机工程与设计, 2011, 32(10): 3570-3573.  
WANG A M, GE W Y, ZHAO Z, et al. Research on computer matching of inscriptions on tortoise fragments[J]. Computer Engineering and Design, 2011, 32(10): 3570-3573. (in Chinese)
- [3] 刘影. 宾组牛胛骨新缀四组[J]. 故宫博物院院刊, 2011, 54(1): 22-27, 156.  
LIU Y. Four sets of newly restored binzu bovine scapula inscription fragments[J]. Palace Museum Journal, 2011, 54(1): 22-27, 156. (in Chinese)
- [4] 齐航福. 甲骨新缀五组[J]. 故宫博物院院刊, 2011, 54(1): 14-21, 156.  
QI H F. A study of five sets of newly restored oraclebone inscription fragments[J]. Palace Museum Journal, 2011, 54(1): 14-21, 156. (in Chinese)
- [5] 王爱民, 刘国英, 葛文英, 等. 甲骨文计算机辅助缀合系统设计[J]. 计算机工程与应用, 2010, 46(21): 59-62.  
WANG A M, LIU G Y, GE W Y, et al. System designation for computer aided rejoining of tortoise shells with inscriptions based on contour matching[J]. Computer Engineering and Applications, 2010, 46(21): 59-62. (in Chinese)
- [6] 孙亚冰. 甲骨缀合五则[J]. 南方文物, 2015, 54(3): 107-108, 82.  
SUN Y B. Five cases of oracle bones splicing[J]. Relics from South, 2015, 54(3): 107-108, 82. (in Chinese)
- [7] 王爱民, 葛彦强, 刘国英, 等. 计算机辅助甲骨文缀合关键技术研究[J]. 计算机测量与控制, 2010, 18(7): 1612-1614.  
WANG A M, GE Y Q, LIU G Y, et al. Research on key technologies of computer aided rejoining of bones/tortoise shells with inscriptions[J]. Computer Measurement & Control, 2010, 18(7): 1612-1614. (in Chinese)
- [8] 王爱民, 钟珞, 葛彦强, 等. 甲骨碎片智能缀合关键技术研究[J]. 武汉理工大学学报, 2010, 32(20): 194-199.  
WANG A M, ZHONG L, GE Y Q, et al. Research on key technologies of the computer aided rejoining of the bones/tortoise shells with inscriptions[J]. Journal of Wuhan University of Technology, 2010, 32(20): 194-199. (in Chinese)
- [9] 王爱民, 葛彦强, 刘国英, 等. 甲骨文计算机辅助缀合技术研究[J]. 中国科技信息, 2010, 22(4): 43-46.  
WANG A M, GE Y Q, LIU G Y, et al. The system designation for the computer aided rejoining of the tortoise shells with inscriptions based on contour matching[J]. China Science and Technology Information, 2010, 22(4): 43-46. (in Chinese)
- [10] 顾绍通. 甲骨文数字化处理研究述评[J]. 西华大学学报(自然科学版), 2010, 29(5): 38-42, 48.  
GU S T. Review on digitization processing of jiaguwen [J]. Journal of Xihua University (Natural Science Edition), 2010, 29(5): 38-42, 48. (in Chinese)
- [11] 张长青, 王爱民. 一种计算机辅助甲骨文拓片缀合方法[J]. 电子设计工程, 2012, 20(17): 1-3.  
ZHANG C Q, WANG A M. Method for computer aided rejoining of bones/tortoise shells rubbing[J]. Electronic Design Engineering, 2012, 20(17): 1-3. (in Chinese)
- [12] RAKTHANMANON T, CAMPANA B, MUEEN A, et al. Searching and mining trillions of time series subsequences under dynamic time warping[C]//Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2012: 262-270.
- [13] MUEEN A, KEOGH E. Extracting optimal performance from dynamic time warping[C]//Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2016: 2129-2130.
- [14] KEOGH E J, PAZZANI M J. Scaling up dynamic time

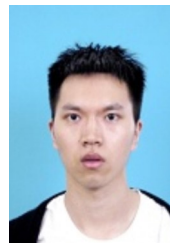
- warping for data mining applications[C]//Proceedings of the sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2000: 285-289.
- [15] JEONG Y S, JEONG M K, OMITAOMU O A. Weighted dynamic time warping for time series classification[J]. Pattern Recognition, 2011, 44(9): 2231-2240.
- [16] RATANAMAHATANA C A, KEOGH E. Three myths about dynamic time warping data mining[C]//2005 SIAM International Conference on Data Mining. Philadelphia: Society for Industrial and Applied Mathematics, 2005: 506-510.
- [17] ZHAO J, ITTI L. shapeDTW: Shape dynamic time warping[J]. Pattern Recognition, 2018, 74: 171-184.
- [18] CAI X, XU T, YI J, et al. DtwNet: A dynamic time warping network[J]. Advances in Neural Information Processing Systems, 2019, 11636-11646.
- [19] SILVA D F, BATISTA G E. Speeding up all-pairwise dynamic time warping matrix calculation[C]//2016 SIAM International Conference on Data Mining. University City, Philadelphia: Society for Industrial and Applied Mathematics, 2016: 837-845.
- [20] ZHANG Z, TAVENARD R, BAILLY A, et al. Dynamic time warping under limited warping path length[J]. Information Sciences, 2017, 393: 91-107.
- [21] JAIN B J. Making the dynamic time warping distance-warping variant[J]. Pattern Recognition, 2019, 94:35-52.
- [22] 杨一鸣, 潘嵘, 潘嘉林, 等. 时间序列分类问题的算法比较[J]. 计算机学报, 2007, 30(8): 1259-1266.  
YANG Y M, PAN R, PAN J L, et al. A comparative study on time series classification[J]. Chinese Journal of Computers, 2007, 30(8): 1259-1266. (in Chinese)
- [23] 孙冬璞, 曲丽. 时间序列特征表示与相似性度量研究综述[J]. 计算机科学与探索, 2021, 15(2): 195-205.  
SUN D P, QU L. Survey on feature representation and similarity measurement of time series[J]. Journal of Frontiers of Computer Science and Technology, 2021, 15(2): 195-205. (in Chinese)
- [24] 周宁南, 张孝, 刘城山, 等. 基于动态时间规整的时序数据相似连接[J]. 计算机学报, 2018, 41(8): 1798-1813.  
ZHOU N N, ZHANG X, LIU C S, et al. Similarity join on time series under dynamic time warping[J]. Chinese Journal of Computers, 2018, 41(8): 1798-1813. (in Chinese)
- [25] 李永健. 基于 DTW 和 HMM 的语音识别算法仿真及软件设计[D]. 哈尔滨: 哈尔滨工程大学, 2009.
- LI Y J. Speech Recognition Algorithm Simulation and Software Design Based on DTW and HMM[D]. Harbin: Harbin Engineering University, 2009. (in Chinese)
- [26] 周瑜, 刘俊涛, 白翔. 形状匹配方法研究与展望[J]. 自动化学报, 2012, 38(6): 889-910.  
ZHOU Y, LIU J T, BAI X. Research and perspective on shape matching[J]. Acta Automatica Sinica, 2012, 38(6): 889-910. (in Chinese)
- [27] SALVADOR S, CHAN P. Toward accurate dynamic time warping in linear time and space[J]. Intelligent Data Analysis, 2007, 11(5): 561-580.
- [28] CORMEN T H, LEISERSON C E, RIVEST R L, et al. Introduction to Algorithms[M]. Cambridge: MIT Press, 2022.
- [29] ARONOV B, HAR-PELED S, KNAUER C, et al. Fréchet distance for curves, revisited[C]//European Symposium on Algorithms. Berlin: Springer, 2006: 52-63.
- [30] BELONGIE S, MALIK J, PUZICHA J. Shape matching and object recognition using shape contexts[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2002, 24(4): 509-522.
- [31] RADENOVIC F, TOLIAS G, CHUM O. Deep shape matching[C]//European Conference on Computer Vision. Berlin: Springer, 2018: 774-791.

#### 作者简介



张重生 男, 1982年生, 河南南阳人. 现为河南大学计算机与信息工程学院教授、博士生导师. 主要研究方向为长尾学习与不平衡学习、汉字识别和古文字计算.

E-mail: eszhang@henu.edu.cn



王斌(通讯作者) 男, 1997年生, 河南安阳人. 河南大学计算机与信息工程学院硕士. 主要研究方向为机器学习.

E-mail: bin.wang@henu.edu.cn