

# 目录

<b>第一章 双分辨率词云的设计目标</b>	<b>1</b>
1.1 任务 . . . . .	1
1.2 设计准则 . . . . .	1
1.3 设计选择 . . . . .	2
<b>第二章 双分辨率词云可视化</b>	<b>5</b>
2.1 数据 . . . . .	5
2.2 布局算法 . . . . .	6
2.2.1 基础布局 . . . . .	6
2.2.2 子文本布局 . . . . .	6
2.2.3 算法加速 . . . . .	7
2.3 混合图像 . . . . .	8
2.3.1 词云融合 . . . . .	8
2.3.2 颜色编码 . . . . .	8
<b>第三章 双分辨率词云创制系统</b>	<b>10</b>
3.1 设计要求 . . . . .	10
3.2 系统概览 . . . . .	11
3.2.1 界面 . . . . .	11
3.2.2 交互逻辑 . . . . .	11
3.3 前后端交互 . . . . .	11
3.3.1 数据上传 . . . . .	11
3.3.2 数据处理 . . . . .	12
3.3.3 参数推荐 . . . . .	13
3.4 系统实现 . . . . .	13
<b>参考文献</b>	<b>14</b>
<b>本科期间的主要工作和成果</b>	<b>16</b>
<b>北京大学学位论文原创性声明和使用授权说明</b>	<b>17</b>



## 第一章 双分辨率词云的设计目标

为了合理提出双分辨率词云及其自定义创制系统的视觉设计，本章节将讨论双分辨率词云所关注的任务与使用情境，结合可视化中相关的重要设计准则，推导出设计要求，以此在众多首要的设计选择中明确具体的方案，明确设计的方向。

### 1.1 任务

双分辨率词云是一种通用的文本可视化方法，其核心是延拓现有词云，提供对上下文的支持，在静态图像中额外编码关键词对应的语境。我们认为，无论是出于探索式数据分析的需求还是出于可视化具有关联的文本的需要，双分辨率词云都能发挥一定的作用。此外，无论是词云还是多分辨率图像，他们都能够娱乐观众，兼具艺术性和趣味性。而脱胎于此的双分辨率词云亦应继承这一良好特性。

- T1. 探索式数据分析**—尽管词云不涉及复杂的自然语言处理技术，仅凭借词频粗略地概括文本内容，但它很适用探索式数据分析的情境。当分析者对数据一无所知，尚未建立任何猜想或假设，关键词这种回归原始数据的底层信息能够使分析者快速构建对数据自身的理解，了解文本中的主要话题。而双分辨率词云则是进一步提供了具体的语境，使得分析者可通过递进的方式进行探索，更为准确地理解数据。
- T2. 层次化文本数据展示**—双分辨率词云的形式除了能够展示关键词及其上下文，也可用于展示其他具有二重层次关联的文本数据。例如外层词云对应主题，内层词云对应主题下的关键词的情况。它能够将抽象的数据处理结果可视化，辅助交流与沟通。
- T3. 寓教于乐**—所谓“一图胜千言”，生动的信息图表（Infographics）往往能激发观众的兴趣，带来别样的浏览体验 [1]。我们希望双分辨率词云具有寓教于乐的功能，在反映数据的同时兼具艺术性与趣味性，让随机浏览的人们沉浸其间，使创作者能更好地向大众传达词云所隐含的信息。

### 1.2 设计准则

除了以上对分析情境及对应任务的考量，我们在双分辨率词云的设计中也需要考虑一些普适的设计准则与认知规律。

#### **P1. 可视化设计标准**

**P1.1 Gestalt理论**—Gestalt理论 [2]描述了人们在感知或理解图像时，根据某些特征将元素聚集所表现出来的规律。通过这种机制，人们得以简化图像的模式、理

解其内涵。其中，与本研究相关的规律有以下三点。

- 相似性聚集：相似的元素会被视为一类。我们习惯于会不自觉将形状或色彩等基本视觉属性接近的元素划归为一类，认为他们具有相同的作用。
- 相邻性聚集：物理位置上靠近的元素会被视为一类。而对于彼此存在巨大空隙的元素，就算他们在外观上具有一定相似性，我们也倾向于认为他们之间无关。即是说，相邻性聚集的作用远远大于相似性聚集。
- 共享空间聚集：处于同一封闭区域内部的元素会被认为同属一类。

**P1.2 自顶而下的探索流程**—可视化领域广泛认可“先概览，通过聚焦和过滤进行筛选，最后提供细节信息”[3]的交互探索流程，大量应用均遵循着这一原则。概览让用户首先对数据产生基本的认识，定位到感兴趣的数据子集，再从细节中完成更为精细的分析任务。

**P1.3 使用更少的视图**。在具有多个视图的可视化中，应尽可能减少所使用的视图数目[4]。尽管利用交互切换视图能帮助用户在数据概览中渐进地定位到关键区域，对更多的细节展开探索，但过多的视图切换伴随着额外的认知开销，会产生学习与交互成本，使用户难以构建对数据的宏观理解（Mental Map），最终影响数据分析的效率。

## P2. 美观标准

**P2.1 空间覆盖率高**—参考Wang等人[5]对词云美观性的评价指标，我们以词云中空白区域像素个数占所有像素个数的比例作为度量。

**P2.2 布局均匀**—在词云的布局中，无关联的文本应均匀地布局在图片上。换言之，词与词之间的空隙是均匀分布的。

**P2.3 字符清晰可辨**—两层词云各不影响。在最理想的情形下，对于特定的距离范围，只有父词云或子词云可见，且同级词云均可清晰辨别。

## 1.3 设计选择

**大屏还是小屏，静态还是交互？** 双分辨率词云的初衷是为基本的词云提供上下文信息，使人们能够结合语境更好地理解数据。以此为出发点，我们其实还有许多设计选择。由于除了客观的人眼感知能力与计算机硬件性能，可视化工具展示数据的能力（即视觉可扩展性[6]）还受到屏幕分辨率、可视系统的交互性、可视化隐喻内涵等因素的影响，我们由此来考虑备选设计。

首先是交互性。在可视分析中，许多系统为用户提供了多种相互关联的视图，让他们通过自行筛选来浏览更多的细节，逐渐发掘数据背后的知识。类似地，我们可以提出这样的一个设计：只显示父词云，再通过动态查询[7]展示相应的子词云信息，如

返回一个列表。结合我们的任务，答案是否定的。这种渐进式的探索需要用户确定感兴趣的区域，再去执行筛选而改变视图，缺乏即时性。对于向多人展示的场景（T2），无法体现数据的全貌，缺乏说服力。而对受众随机浏览的场景（T3），由于交互不是显式存在的，极有可能会被忽略，最终与静态普通词云的效果无异，背离了初衷。

而由于文本的非结构化特性，为其换用图形表征是较为困难的，因此不予考虑。最终，我们选择在大屏上静态地同时显示上下层词云。提高屏幕分辨率后极大遍历了大量信息的同时显示。在现有的任务框架下，针对大屏的静态可视化设计是最佳选择。借助引人入胜的多分辨率技术（T3），其只需要一个视图（P1.3），较易使人理解。人们在远处即可获得对整体的认知，确定感兴趣区域后只需物理上稍微靠近屏幕即可获取更多的细节（P1.2），非常适合多人探索的任务（T1，T2）。

**如何编码字号？** 词云中字的大小与对应的权重（如词频）具有明确的一对一关系，其作用是将权重具有显著差异的词目区分开来。由于词的长度不一，且词与词之间的相对位置不是对齐的，人们很难通过字号来准确地对比两个相似大小的词之权重 [8]。

存在多种权重与字号的对应关系 $f$ ，满足单调递增且值域非负即可，如斜率为正的线性函数、排名函数、开平方函数等。我们认为，何种映射适于区分权重的层次是由数据自身决定的。例如，当存在一个特别大的权重异常值时，取开平方的即可缓和其与其他数据的差异，理论上优于一视同仁的线性函数。为了提高多分辨率词云对不同数据的适应能力，适应探索（T1）与呈现（T2，T3）的需求，我们认为不应局限权重与字号对应的形式，用户应能够自主选择合理的 $f$ 。

根据Isenburg [9]对WILD设备的实验结果（见图 ??右侧），64像素是在3米开外区分远近文字的一个合理阈值。类似地， $f$ 的值域是相对显示屏幕尺寸及分辨率固定的，只需进一步确定 $f$ 的函数类以及数据的定义域，即可构建映射。

**文字方向有何约束？** 早年的标签云多为水平对齐，随着Wordle的诞生，文本任意旋转的词云获得了更为广泛的关注。然而，在我们大屏的设定下，由于在接近屏幕时，浏览较远处会产生一定的视角扭曲，为了使文字更易被识别（P2.3），我们限定词云为全水平布局。这也能简化基于搜索的词云布局算法，加快大屏上词云布局的计算。

**下层词云如何布局？** 以ShapeWordle [5]为代表的形变词云能够让第二级的文本紧凑布局于上级文本字形所天然形成的边界内部。但由于文字本身的空间覆盖率较低（P2.1），凑近时并不便阅读。且这种方法能够涵盖的子词云较少，视觉可扩展性低。而对图像拼接的方法来说，子文本需要在上层形状的限制下稍加变形，难以控制其字号与权重的对应关系，且不易读。

我们选择在父词云所确定的领域中放置子词云，以空间的相邻表示他们之间的关联（P1.1-2），指引自顶而下的探索（P1.2）。同时，我们还应保证子词云尽可能多地利用空白的位置（P2.1），并均匀排布（P2.2）。

**如何使用颜色？** 尽管不同的灰度值足以实现双分辨率词云的基本要求，但我们仍选择使用颜色来为其增效。鲜艳多彩的颜色不仅能引人注目（T2，T3），更能作为一个单独的视觉通道编码数据（T1，T2）：色相可对应于离散型变量，而亮度或饱和度的变化在一定程度上也可对应于连续性变量。在一般的词云中，词的颜色仅用以区分各个短语。当考虑词义 [10, 11] 或情感 [12] 等附加的属性时，色相有时会被用来表示一个类。在多分辨率词云中，类比于大多数词云的做法，我们使用不同的色相随机编码父词云。但对于与父词云具有关联的子词云，我们在父词云颜色的基础上添加扰动，用相似但稍有区分度的颜色来编码子词云，以此保持认知上的关联性（P1.1）。进一步地，我们还能通过调整词的亮度，区分开父词云与子词云，缓解他们对彼此的干扰（P2.3）。

综上，我们基于双分辨率词云的具体任务与一般的设计准则对其生成时的一些关键问题进行了分析，作出了以下设计选择：

- C1. 考虑大屏上的静态可视化。**充分利用大屏的高分辨率特性，为多人浏览场景提供服务，为理解抽象文本尽多地提供线索。
- C2. 保留权重信息，根据数据特点灵活选择映射类型。**为了将字的权重有效分层，双分辨率词云应在保留原始权重值的基础上给予用户调整权重-字号映射的空间。
- C3. 水平布局。**防止大屏伴随的视角扭曲问题影响浏览，并提高算法效率。
- C4. 子词云布局于父词云的邻域内。**在大屏上尽可能多地展示数据，同时保持父子词云在视觉上的关联性。
- C5. 层次化赋色。**以差异较大的色相相区分父词云，在父词云颜色的基础上添加扰动分别子词云。其中，子词云亮度随其与父词云相对位置的改变有所调整。

## 第二章 双分辨率词云可视化

我们从上文讨论的设计选择出发，在基准数据上进行了一些先行的尝试，通过不断迭代明确了最终双分辨率词云的可视化形式以及相应参数。

### 2.1 数据

双分辨率词云的基本数据模型可对应于二分图（见图 2.1），父级文本（如关键词）与子级文本（如上下文）之间存在层级的关系或语义上的关联。图中的每一个结点均具有字符串型文本属性和数值型的权重或词频属性（C2）。为了减少数据存储的冗余，使数据简洁而结构明晰，我们选择使用两个序列作为双分辨率词云的表示，分别对应于父文本和子文本，而关联信息编码在父文本各项中。

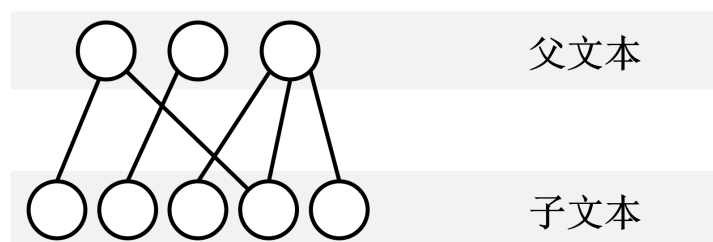


图 2.1 双分辨率词云数据模型。

使用YAML语言，双分辨率词云所对应的数据模型可形式化描述如下。需要注意的是，这里父文本与子文本的权重已归一化。在常见的词云算法中，归一化一般为词频直接除以最大值，亦有取对数或开平方后采用最小-最大归一的方法以降低极端值对整体的影响。

```
1 # 根对象
2 type: object
3 fields:
4   # contexts域：序列类型，每个对象对应一个上下文
5   contexts:
6     type: sequence
7     fields:
8       #context域：子文本，上下文
9       context: { type: string }
```

```

10     #weight域: 子文本权重
11     weight: { type: number, min: 0, max: 1 }
12 # keywords域: 序列类型, 每个对象对应一个关键词
13 keywords:
14     type: sequence
15     fields:
16         # keyword域: 父文本, 关键词
17         keyword: { type: string }
18         # weight域: 关键词权重
19         weight: { type: number, min: 0, max: 1 }
20         # contexts域: 集合类型, 对应的子文本序列号
21         context:
22             type: set
23             fields:
24                 # contextID域: 子文本序列号
25                 contextId: { type: number }

```

Listing 2.1 数据模型

我们选用了《全唐诗》的数据作为基准, 提取词频最高的30词作为父文本, 对应的诗句(逗号或句号为分句符)作为子文本, 并整理为以上数据模型所规定的格式进行尝试。

## 2.2 布局算法

我们依次确定父文本词云与子文本词云的位置。其中, 父词云由已有词云算法给出, 子词云则在对应父元素位置与字形的限制下布局。

### 2.2.1 基础布局

抽象而言, 词云的布局算法就是在给定画布尺寸 $w \times h$ (长 $\times$ 高)与词集合 $\{w_i\}$ 及其对应字号后 $\{s_i\}$ , 确定每个词在画布上的位置 $\{(x_i, y_i)\}$ 。

常见的词云算法在第2.3节已有概述, 它们各有优劣, [伪代码](#)

### 2.2.2 子文本布局

区域分割 [伪代码](#), [描述](#)



分步布局 介绍形态学，图，伪代码，效果

### 2.2.3 算法加速

大屏上的双分辨率词云布局中最主要的问题就是分辨率过大，计算复杂。无论是使用何种算法，其复杂度都与画布尺寸相关，全屏情况即对应屏幕分辨率。由于无法避免检测词语位置是否重叠，各算法需要对图像进行扫描，逐像素确定合适位置；而对于迭代型算法如力导向布局，其计算更是代价高昂。因此，尽管以上的布局策略理论上可行，但在实际针对大屏的计算中，我们仍有一定的优化空间，提升算法效率。

首先，父文本词云的精确计算是不必要的。我们注意到，父文本词云具有相当大的字号，在远处浏览父词云的效果与在近处浏览常规词云的效果类似。换言之，大屏上的父词云实际可视为整体放大的常规词云。因此，我们可以对父词云进行适当的全局放缩，在更小的画布上计算对应小字号词云的布局，再相应地映射到原有的大屏上，做一个简单的近似。

全局等比放缩系数 $k$ 与父文本最小字号 $m$ 有关。在设计选择的讨论中，我们已经确定父文本字号阈值设为64是合适的，而考虑到由于字号为离散整数，过度的放缩会限制词云字号的表示空间，也会产生过度失真的现象（见图 2.2），因此我们设定16像素为缩放后的最小字号，即 $k = 64/16 = 4$ 。

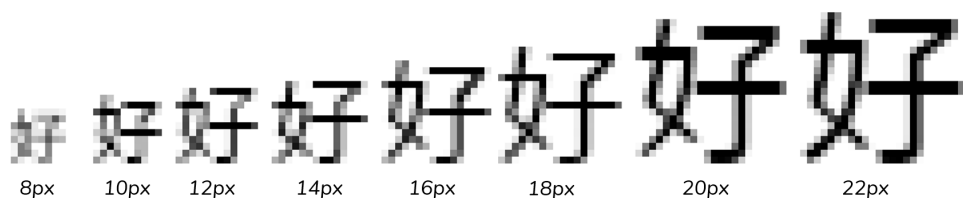


图 2.2 不同分辨率下的文字“好”，从左至右由10px至22px。分辨率越低，字形失真越严重。

例如，对一个由4×3个4K显示屏所组成的大型显示设备，其具有 $w_1 \times h_1 = 15360 \times 6480$ 的分辨率。我们通过缩放，基于 $w_2 \times h_2 = 3840 \times 1620$ 的画布尺寸得到每个词的基础坐标 $\{(x'_i, y'_i)\}$ 。那么，经过一个线性变化，即可得到在原始画布上的布局坐标估计值：

$$\begin{pmatrix} x_i \\ y_i \end{pmatrix} = k \begin{pmatrix} x'_i - w_2/2 \\ y'_i - h_2/2 \end{pmatrix} + \begin{pmatrix} w_1/2 \\ h_1/2 \end{pmatrix}.$$

此外，子词云布局算法可数据并行。由于在分割区域时已保证各子词云位置无重叠，因此各父文本对应的子文本词云的计算是完全独立的，可充分利用大型显示设备的GPU资源等来进行并行化计算，加速双分辨率词云生成的过程。

## 2.3 混合图像

### 2.3.1 词云融合

双分辨率词云的关键在于使人们在不同的距离下关注到不同层次的词云。参考Olivia等人 [13]的多分辨率图像生成算法，我们同样采用先过滤父文本词云的高频分量，再通过透明度叠加的方法混合上下两个图层的方法来得到最终的结果。

由于子文本在双分辨率词云的设计中本就具有较小的字号，人们只有凑近屏幕至适宜阅读距离才能辨别，而滤波会影响其字形，因此我们不对其空间频域进行特殊处理。而对于字号较大的父文本，人们是从远处浏览，无需保留边界的细节即可辨认，且滤波后能够降低其对识别子词云产生的干扰，所以有必要过滤其高频分量。

简记在上一节所述布局算法得到的父文本词云为 $I_1$ ，子文本词云为 $I_2$ 。我们使用高斯滤波器对 $I_1$ 进行卷积。一个标准的高斯函数 $G(x, y)$ 的表达式为

$$G(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}},$$

其中 $\sigma$ 表示其标准差。一个高斯卷积核就是以核中心为原点的经过和归一化的离散高斯函数。设 $I_1$ 的最小字号为 $s_{min} = \min_i s_i$ ，则我们限定高斯核 $Ker(G)$ 的大小为 $k^2 = \text{ceil}(\frac{s_{min}}{25})^2$ 。

$$I'_1 = I_1 \otimes Ker_k(G).$$

$$I'_2 = I_2.$$

$k$ 分母上的常数25是我们通过类似图 2.2的实验确定的，在常见的一些字体中，当分辨率 $r$ 逐渐变大，对应字符比特图上的笔触宽度约为 $\text{ceil}(\frac{r}{25})$ 。我们在全局使用与最小词云笔触大小一致的卷积核以防止字形剧变，影响识别。

经过 $\alpha \in [0, 1] \cap \mathbb{R}$ 的透明度叠加，最终得到的图像 $I$ 可表示为

$$I = \alpha I'_1 + (1 - \alpha) I'_2.$$

一个流程图

### 2.3.2 颜色编码

彩色图片多以RGB（分别对应红、绿、蓝）三个色彩通道来编码一个像素所对应的颜色。

尽管RGB编码形式有利于计算机上的数字图像处理，但它人的视觉感知规律大相

径庭。人是通过视网膜上分布的视杆细胞与视锥细胞来感知光线的，分别对应低亮度与高亮度的情况。Opponent-Process理论 [14]认为，大脑将感光细胞做感知到的光波处理为了三个通道——亮度、红绿与蓝黄通道。以此为理论指导，国际照明委员会定义了LAB色彩空间以模拟人的感知，经笛卡尔坐标系转换为圆柱坐标系后得到LCH色彩空间（见图 2.3）。

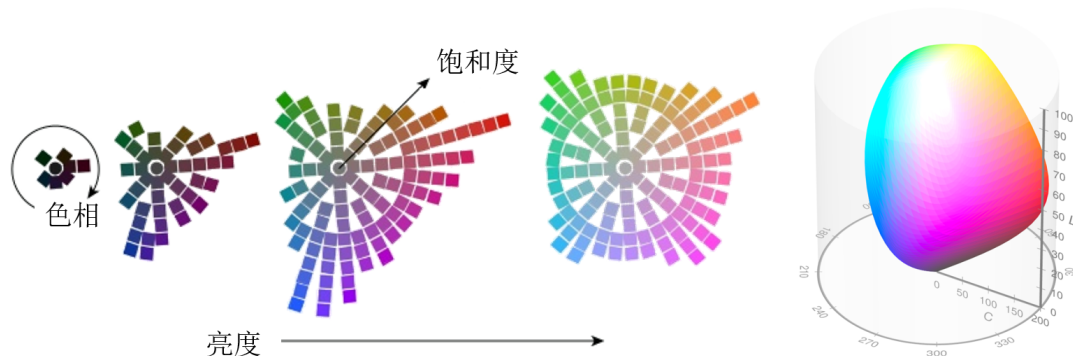


图 2.3 色彩空间。左：不同亮度下人眼所能感知到的颜色示意（图源：NASA Earth Observatory<sup>①</sup>）。右：LCH色彩空间（图源：维基百科<sup>②</sup>）。

② Subtleties of Color. Robert Simon. <https://earthobservatory.nasa.gov/blogs/elegantfigures/2013/08/05/subtleties-of-color-part-1-of-6/>

② HCL color space. [https://en.wikipedia.org/wiki/HCL\\_color\\_space](https://en.wikipedia.org/wiki/HCL_color_space)

## 第三章 双分辨率词云创制系统

双分辨率词云是一种任务驱动型的基础可视化，可应用于相当广泛的文本数据。尽管它可以在上一章算法的基础上被独立开发实现，但由于它是一种新颖的可视化方法，且所涉及的参数较多，新接触的人群难以理解参数与最终效果的关联，且尚无完备的优化准则来指导机器自动生成，因此不利于创作者与分析人员的了解与利用。为此，我们设计开发了一套基于用户-服务器模型的交互创制系统，词境，以支持用户对双分辨率词云的便捷创制。

### 3.1 设计要求

基于在设计双分辨率词云可视化的迭代过程中所积累的一些经验以及对现有的在线词云创建系统的参考，我们总结出了以下几项对双分辨率词云创制系统的设计要求，以方便用户快速生成双分辨率词云，并在一定程度上控制最终的效果。

- D1. 提供灵活的数据处理功能：**由于双分辨率词云适用于探索式分析任务（T1）和数据分享任务（T2），因此创制系统也应相应地对数据格式有更灵活的支持，为具有不同需求的用户提供服务。
  - D1.1 纯文本语料库—**对毫无先验知识的纯文本，用户可直接上传，系统在简单处理后可可视化其高频词及对应语句或文短以提供具有语境的概览，帮助用户找到数据中感兴趣的内容，挖掘数据背后的洞见。
  - D1.2 格式化文本—**对数据已具备一定了解的用户亦可指定双层词云各自的具体内容，构建映射关系，展示相应内容，促进交流。在这种情境下，系统应不再局限于简单的文本挖掘，支持用户自定义所需要可视化的内容。
- D2. 支持对效果的调整：**在数据分享任务（T2）和寓教于乐任务（T3）中，可视化的最终效果需要符合设计者一定的心理预期，支持设计者对布局及格式的调整。
  - D2.1 参数推荐—**过于灵活的参数设置容易使可视化经验较少的用户茫然无措，而本研究所提出的双分辨率词云更是一种全新的方法，我们尚未充分地认识到参数与效果之间的对应关系。因此，我们需要在交互的过程中为用户提供更多的指引，基于固定的在不影响用户与系统交互的效率的前提下，尽可能地缩小用户需要探索的参数空间。
  - D2.2 反馈效果—**及时反馈参数所对应的效果，避免让用户通过不断生成双分辨率词云的低效方式来调试。

**D2.3 自定义**—用户可具体指定各词语的可视化参数，如父文本词云（或与子文本词云）的布局位置或赋色方案。

## 3.2 系统概览

词境是一个轻型的网站，由四个页面组成：主页、创制页面、画廊及说明页，采用前后端分离的客户-服务器模式。

### 3.2.1 界面

### 3.2.2 交互逻辑

## 3.3 前后端交互

### 3.3.1 数据上传

本节简要介绍该系统所支持的数据格式与布局参数接口。在下一章里，我们将通过几个具体的例子，说明该设定如何能够方便用户运用双分辨率词云可视化。

**格式化文本** 在前一章的数据模型设计下，词境以json作为双分辨率词云的底层数据格式，并接收用户上传的符合规范的json数据（格式规范如下所示），以适应于更广泛的数据类型（D1.2）。这是一种独立于程序的数据交换语言，具有以下优点。

- **易读性强**—它的层次结构清晰而简洁，易于理解。
- **普及性高**—json是ECMAScript®语言规范的标准之一，受几乎所有编程语言的支持，用户可在简单的数据预处理后生成相应文件。
- **适合网络上的数据交换**—在客户端与服务器端数据传输的过程中针对json型数据存在特定的压缩技术，能够节省带宽。

在此基础上，我们给予用户定制可视化的空间，允许他们指定父级词云的位置与颜色（D2.3）。

```
1 {  
2   "keyword_dict": {  
3     word: {  
4       weight: weight_of_word,  
5       contexts: [  
6         context_id  
7       ],
```

```
8     position: [x, y],
9     color: color_string
10 },
11 "context_list": [
12     context_id: {
13         "context": context_string,
14         "weight": weight_of_context
15     }
16 ]
17 }
```

Listing 3.1 json上传数据格式规范，父级词云的位置与颜色域可为空

**纯文本** 在常见的词云创制网站，如Wordle<sup>①</sup>和Word Art<sup>②</sup>中，用户通过复制粘贴文本或指定文字来源网页超链接的方式来指定可视化的内容，限定使用UTF-8编码。注意到它们对于电子书、预印本等常见语料库的支持不足，我们还在系统内额外增添了pdf和txt格式的文件处理机制（D1.1），将他们转化为字符串。

## 布局参数

### 3.3.2 数据处理

对于用户所上传的无任何先验知识的纯文本，系统默认将词频高的关键词作为父词云，其所在句子作为子词云。为了进一步确定用户所关心的细节程度，我们需要用户输入句子的分隔符，默认为句号、分号、感叹号与问号。

中文、日文、泰文等本文的一大特色在于，它们的句段中短语或词组不存在天然的分隔符。而由于语言的歧义性、词库有限等诸多难题，如何准确分词至今仍是热门的研究方向。为此，系统使用了性能较为良好的“结巴分词”<sup>③</sup>进行了分词操作以提取关键词。

---

① <http://www.wordle.net>

② <https://wordart.com>

③ <https://github.com/fxsjy/jieba>

### 3.3.3 参数推荐

## 3.4 系统实现

本研究所实现的双分辨率词云创制系统的HTML5网站在前端部分使用了轻量级的Vue框架生成动态页面，以JavaScript语言协调用户与系统的交互，包括数据上传、参数输入、输入完整性检验与推荐，以及生成图像预览与下载。而后端部分则是基于Flask框架响应前端请求，对于自定义数据进行相应的处理，并根据参数返回生成的词云。后端所用程序语言为Python，单独开发了双分辨率词云的类库，提供函数接口，在词云布局的迭代算法中使用了Cython以进行加速。

[采访系统使用感受。](#)

## 参考文献

- [1] Tamara Munzner. *Visualization Analysis and Design*. A K Peters, **2014**. <http://www.cs.ubc.ca/%5C%7Etm/vadbook/>.
- [2] Johan Wagemans, James H. Elder, Michael Kubovy *et al.* “A century of Gestalt psychology in visual perception: I. Perceptual grouping and figure-ground organization.” *Psychological Bulletin*, **2012**, 138(6): 1172–1217. <https://doi.org/10.1037/a0029333>.
- [3] Ben Shneiderman. “*The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations*”. In: *Proceedings of the 1996 IEEE Symposium on Visual Languages, Boulder, Colorado, USA, September 3-6, 1996*. IEEE Computer Society, **1996**: 336–343. <https://doi.org/10.1109/VL.1996.545307>.
- [4] Michelle Q. Wang Baldonado, Allison Woodruff and Allan Kuchinsky. “*Guidelines for Using Multiple Views in Information Visualization*”. In: *Proceedings of the Working Conference on Advanced Visual Interfaces*. Palermo, Italy: Association for Computing Machinery, **2000**: 110–119. <https://doi.org/10.1145/345513.345271>.
- [5] Yunhai Wang, Bongshin Lee, Xiaowei Chu *et al.* “*ShapeWordle: Tailoring Wordles using Shape-aware Archimedean Spirals*”. *IEEE Transactions on Visualization and Computer Graphics*, **2020**, 26(1): 991–1000. <https://doi.org/10.1109/tvcg.2019.2934783>.
- [6] Stephen G Eick and Alan F Karr. “*Visual Scalability*”. *Journal of Computational and Graphical Statistics*, **2002**, 11(1): 22–43. <https://doi.org/10.1198/106186002317375604>.
- [7] Ben Shneiderman. “*Dynamic Queries for Visual Information Seeking*”. *IEEE Software*, **1994**, 11(6): 70–77. <https://doi.org/10.1109/52.329404>.
- [8] Johann Schrammel, Michael Leitner and Manfred Tscheligi. “*Semantically Structured Tag Clouds: An Empirical Evaluation of Clustered Presentation Approaches*”. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Boston, MA, USA: Association for Computing Machinery, **2009**: 2037–2040. <https://doi.org/10.1145/1518701.1519010>.
- [9] Petra Isenberg, Pierre Dragicevic, Wesley Willett *et al.* “*Hybrid-image visualization for large viewing environments*”. *IEEE Transactions on Visualization and Computer Graphics*, **2013**, 19(12): 2346–2355. <https://doi.org/10.1109/TVCG.2013.163>.
- [10] Lukas Barth, Stephen G. Kobourov and Sergey Pupyrev. “*Experimental Comparison of Semantic Word Clouds*”. In: *Proceedings of the 13th International Symposium on Experimental Algorithms - Volume 8504*. Springer-Verlag, **2014**: 247–258. [https://doi.org/10.1007/978-3-319-07959-2\\_21](https://doi.org/10.1007/978-3-319-07959-2_21).
- [11] Marti Hearst, Emily Pedersen, Lekha Priya Patil *et al.* “*An Evaluation of Semantically Grouped Word Cloud Designs*”. *IEEE Transactions of Visualization and Computer Graphics*, **2019**. <https://doi.org/10.1109/TVCG.2019.2904683>.



- [12] Tugba Kulahcioglu and Gerard de Melo. “Paralinguistic recommendations for affective word clouds”. In: *Proceedings of the 24th International Conference on Intelligent User Interfaces*. ACM, **2019**. <https://doi.org/10.1145/3301275.3302327>.
- [13] Aude Oliva, Antonio Torralba and Philippe G. Schyns. “Hybrid Images”. *ACM Transactions of Graphics*, **2006**, 25(3): 527–532. <https://doi.org/10.1145/1141911.1141919>.
- [14] Leo M Hurvich and Dorothea Jameson. “An opponent-process theory of color vision.” *Psychological Review*, **1957**, 64(6): 384–404. <https://doi.org/10.1037/h0041403>.

## 本科期间的主要工作和成果

### 已接收论文

1. **Liwenhan Xie**, James O'Donnell, Benjamin Bach, and Jean-Daniel Fekete. Iterative Time-Series of Measures for Exploring Dynamic Networks. *International Conference on Advanced Visual Interfaces*, Island of Ischia, Italy, 2020.
2. Can Liu, **Liwenhan Xie**, Yun Han, Datong Wei and Xiaoru Yuan. AutoCaption: An Approach to Generate Natural Language Description from Visualization Automatically. *IEEE Pacific Visualization Symposium (Notes)*, Tianjin, China, 2020.

### 获奖作品/海报

1. Shuai Chen, Sihang Li, **Liwenhan Xie**, Yi Zhong, Yun Han, and Xiaoru Yuan. EarthquakeAware: Visual Analytics for Understanding Human Impacts of Earthquakes from Social Media Data. *IEEE Symposium on Visual Analytics Science and Technology (VAST Challenge)*, Vancouver, Canada, 2019. **Honorable Mention for Support for Analysis through Annotation and Context Award.**
2. Can Liu, **Liwenhan Xie**, Yun Han, Datong Wei and Xiaoru Yuan. Automatic Caption Generation for SVG Charts. *IEEE Pacific Visualization Symposium*, Bangkok, Thailand, 2019. **Honorable Mention for Best Poster Award.**
3. Qi Ma, Chuangming Huang, **Liwenhan Xie**, Zhiyi Yin ,and Xiaoru Yuan. Visual Analysis for Subgroups in a Dynamic Network. *IEEE Symposium on Visual Analytics Science and Technology (VAST Challenge)*, Berlin, Germany, 2019. **Insights Generated Through the Use of a Custom Tool Award.**
4. Lijing Lin, **Liwenhan Xie**, Zhuo Zhang, and Xiaoru Yuan. Visualizing Dynamic Networks of Long Sequences with Pixel Matrix Array. *China Visualization and Visual Analytics Conference*, Sichuan, China, 2018. **Honorable Mention for Best Poster Award.**

## 北京大学学位论文原创性声明和使用授权说明

### 原创性声明

本人郑重声明：所呈交的学位论文，是本人在导师的指导下，独立进行研究工作所取得的成果。除文中已经注明引用的内容外，本论文不含任何其他个人或集体已经发表或撰写过的作品或成果。对本文的研究做出重要贡献的个人和集体，均已在文中以明确方式标明。本声明的法律结果由本人承担。

论文作者签名：                    日期：    年    月    日

### 学位论文使用授权说明

（必须装订在提交学校图书馆的印刷本）

本人完全了解北京大学关于收集、保存、使用学位论文的规定，即：

- 按照学校要求提交学位论文的印刷本和电子版本；
- 学校有权保存学位论文的印刷本和电子版，并提供目录检索与阅览服务，在校园网上提供服务；
- 学校可以采用影印、缩印、数字化或其它复制手段保存论文；
- 因某种特殊原因须要延迟发布学位论文电子版，授权学校在□一年/□两年/□三年以后在校园网上全文发布。

（保密论文在解密后遵守此规定）

论文作者签名：                    导师签名：                    日期：    年    月    日