

Lab4实验报告

实验题目：最短路径中文文本分词

姓名：侯新铭

学号：2021201651

1. 需求分析

中文文本没有自然分界符，对文本语义的解析需要进行（相较于英文等语言）额外的分词`segmentation`操作。本次lab目标基于词库，实现一个输入中文文本，返回成功分词后文本的基础软件。

(1) 输入

在项目`segmentation.exe`程序的输入框中输入待分词中文文本。

可支持文本长度：经文本长度测试，单次切词最多可支持约1万字文本

运行时间：几乎可忽略不计

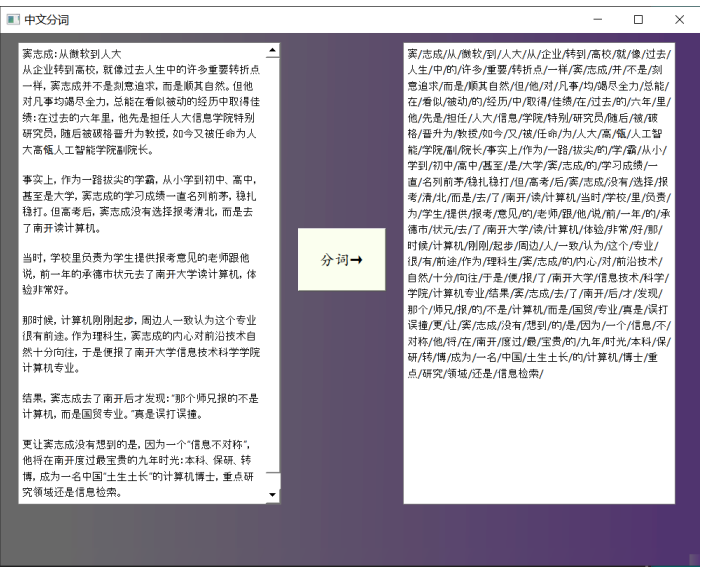
(2) 输出

在交互界面的输出框显示分词后的字符串，格式为以"/"间隔的纯汉字文本。

(3) 功能实现

我开发了一个基础的软件程序，在其可视化界面中，用户可交互式输入中文文本，点击分词按钮后即可立即看到分词后结果；后续用户直接更改其输入的文本内容，实现实时分词。

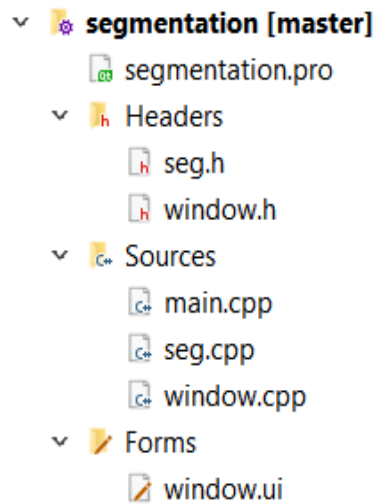
(4) 使用样例



2. 概要设计

项目架构

①使用Qt Creator开发项目主体，项目文件包括：



②头文件调用关系为:

```
|-- main.cpp
    |-- window.h
        |-- seg.h

|-- window.cpp
    |-- window.h

|-- seg.cpp
    |-- seg.h
```

③头文件的主体分别为：

seg.h

```
#ifndef SEG_H
#define SEG_H
#include <queue>
#include <string>
#include <array>
#include <set>
using namespace std;
#define LEN 640
#define INF 0X7FFFFFFF
class Seg
{
public:
    // 载入dict.txt到set类型的dict变量
    bool loadDict(const string &location);
    string cut(string &s);
```

```
private:
    set<string> dict;
};

#endif
```

window.h

```
#ifndef WINDOW_H
#define WINDOW_H

#include <QMainWindow>
#include "seg.h"

QT_BEGIN_NAMESPACE
namespace Ui { class window; }
QT_END_NAMESPACE

class window : public QMainWindow
{
    Q_OBJECT

public:
    window(QWidget *parent = nullptr);
    ~window();

private slots:
    void on_pushButton_clicked();

private:
    Ui::window *ui;
    Seg seg;
};

#endif
```

3. 详细设计

(1) 分词功能实现

④词典加载

```
#include "seg.h"
#include <fstream>
#include <codecvt>
#include <locale>

// 载入dict.txt到set类型的dict变量
```

```

bool Seg::loadDict(const string &location)
{
    ifstream fin(location); // 通过ifstream流读取文件
    if (!fin.is_open())
    {
        return false;
    }

    string line;           // 将文件逐行读取到字符串line中, 截取第一部分
    while (!fin.eof()) // 读取到文件末尾的EOF前一直执行while循环
    {
        getline(fin, line);
        int end = line.find_first_of(' '); // end对应第一列后的空格的索引下标, 也即单词字符串的长度
        if (end != -1) // find_first_of()函数返回值不为-1, 即找到了所给' '字符
        {
            dict.insert(line.substr(0, end)); // 将该词语子串插入到类内的private变量dict中
        }
    }
    return true;
}

```

②分词函数

采用最短路径匹配分词算法实现, 主体循环遍历文本全部字符, 其内包含如下3部分:

- 借助分隔符切分出由中文字符构成的单句
- 基于词典构建单句对应的有向无环图
- 执行dijkstra算法进行分词
- 将该单句切分出的词逆序添加到路径向量中

遍历完全部单句, 主体循环结束, 得到了由全部分词构成的路径向量, 添加“/”作为分隔符, 转化为分词后字符串, 即为所求。

```

// 使用dijkstra算法获得最短路径
string Seg::cut(string &s)
{
    int startPos = 0;
    int fullLen = s.size();
    vector<string> path; // 使用vector记录最终所求的最短路径, 方便进行插入、倒置等操作
    int count = 0;      // 用来记录上一句插词结束后path中的词数, 作为下一句向path中插词的位置标记

    while (startPos < fullLen)
    {
        // 判断中文字符方式: 基于中文字符由3字节构成, 转化为unsigned int必然大于0x7f
        while (!((unsigned int)s[startPos] > 0x7f)) // 找到首中文字符的位置作为

```

sentence的起始

```
{
    startPos++;
    continue;
}
if(startPos>=fullLen) break;
int endPos = startPos;
do
{
    endPos += 3; // 注意到
    // 一个中文字符在utf-8中占3个字节，故以3为步长
} while (endPos < fullLen and ((unsigned int)s[endPos] > 0x7f)); // 找到连续中文字符末位置作为sentence的末尾的下一位
string sentence = s.substr(startPos, endPos - startPos);

int num = sentence.length() / 3; // 恰好为当前sentence包含的词数
array<int, LEN> g;
g.fill(-1);
array<array<int, LEN>, LEN> graph;
graph.fill(g); // 定义2dim array graph,索引到的数值记录sentence各位置之间的可达性。每个位置均初始化为-1, 表示不可达

for (int i = 0; i < num; i++)
{
    graph[i][i + 1] = 2; // 每个字符和下一个字符显然可连通，距离为2，即闭区间跨越2个词可达到
}

for (int i = 0; i <= num - 2; i++)
{
    for (int j = 2; j <= num - i and j <= 12; j++) // dict最长词长度为12字符，作为查找上限
    {
        string checkStr = sentence.substr(i * 3, j * 3);
        if (dict.count(checkStr))
        {
            graph[i][i + j] = 1; // 更新graph, 表示i节点和i+j节点可处于同一个词内，距离为1
        }
    }
}
// 下述部分执行dijkstra算法过程，对当前的sentence进行切词
// 初始化
array<int, LEN> d; // 记录dijkstra算法执行到当前时刻，各节点到初始节点的距离
d.fill(INF);
d[0] = 0;
array<int, LEN> preNum; // 记录各节点当前所得的最小路径的前驱节点序号
preNum.fill(-1);
array<int, LEN> used; // 记录是否已作为最短路径节点使用过
used.fill(0);
for (int i = 1; i <= num; i++) // 初始化与初始节点直接相连的各节点
{
    if (graph[0][i] > 0)
    {
```

```

        d[i] = graph[0][i];
        preNum[i] = 0;
    }
}
// dijkstra算法主体
for (int i = 1; i <= num; i++) // 遍历寻找到未被使用过的节点中距初始节点的最短
距离, 该节点序号记为k
{
    int min = INF;
    int k = 0;
    if (!used[i] and d[i] < min)
    {
        min = d[i];
        k = i;
    }
    used[k] = 1; // k节点被使用, 更新used

    for (int j = 1; j <= num; j++)
    {
        if (graph[k][j] > 0 and graph[k][j] + min < d[j]) // 借助k节点可以
构造出到节点j到初始节点的更短距离, 更新d和preNum
        {
            d[j] = graph[k][j] + d[k];
            preNum[j] = k;
        }
    }
}

while (num != 0) // 从当前centence最后字符处往前不断找前驱节点, 过程中把词
edges构建为路径vector path
{
    string word = sentence.substr(preNum[num] * 3, (num - preNum[num]) *
3); // 截取以当前位置为末节点的词edge
    wstring_convert<codecvt_utf8<wchar_t>> cov;
    // 欲借助c++的iswpunct()和iswspace()来判断是否为中文标点, 需先将string转化为wchar_t
    wchar_t ch = cov.from_bytes(word)[0];
    if (!iswpunct(ch) && !iswspace(ch)) // 判断为无(非)中文标点的word, 插
入path中固定位置count处, 自然地实现了逆序插入
    {
        path.insert(path.begin() + count, word);
    }
    num = preNum[num];
}
count = path.size(); // 完成当前centence的分词后, 更新位置标记count
startPos = endPos; // 更新startPos, 向后推进
}
string output;
for (int i = 0; i < count; i++) // 返回path中各词用"/"连接后的字符串
{
    output.append(path[i] + "/");
}
return output;
}

```

(2) GUI界面开发

①使用Qt Creator创建Qt Widgets Application类型项目，生成主函数main.cpp为：

```
#include "window.h"

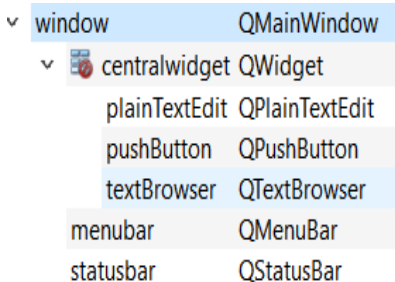
#include <QApplication>

// 使用Qt Creator创建project后自动生成main代码，执行后，运行应用程序
int main(int argc, char *argv[])
{
    QApplication a(argc, argv);
    window w;
    w.show();
    return a.exec();
}
```

②在项目中新建Q MainWindow类型文件组，命名为window，其中包括：

```
-- window.h
-- window.cpp
-- window.ui
```

③进而使用Qt Designer设计GUI界面，主要包括输入框、输出框、push按钮三个组件，如下图所示：



④在window.h中的class window中实例化已经写好的class Seg以调用loadDict和cut函数，最后实现qt中的接口槽函数on_pushButton_clicked()：

```
void window::on_pushButton_clicked()
{
    string s = ui->plainTextEdit->toPlainText().toStdString();
    string output = seg.cut(s);
    ui->textBrowser->setText(QString::fromStdString(output));
}
```

4. 调试分析

时间复杂度：

- dict通过哈希表形式的set定义, 调用check函数的时间复杂度为 $O(1)$;
- 对于有向图 $G=(V, E)$ 记 n 为节点数, dijkstra算法在分词最短路径问题中, 搜索 $d[]$ 中最小距离需 $O(n)$, 内层循环构建新edge, 更新参数的循环也需 $O(n)$, 算法总体时间复杂度为是 $O(n^2)$ 。

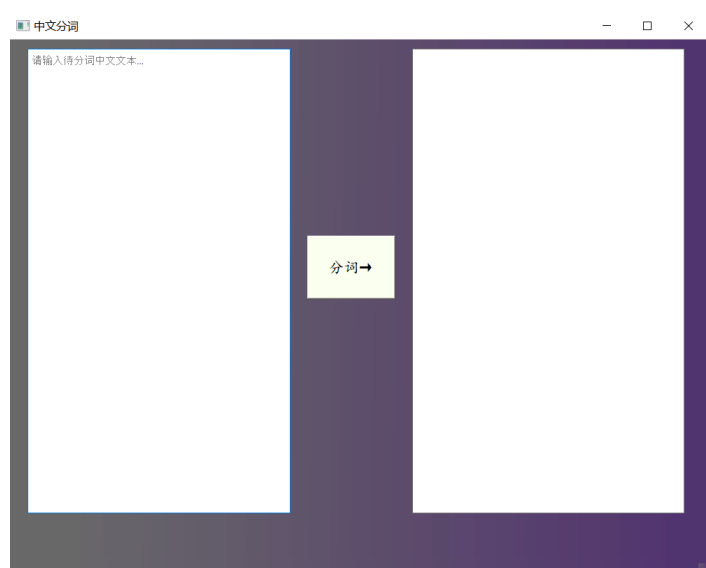
其他通过调试解决了的问题已在代码中详细注释。

5. 用户手册

- 项目在Window系统上开发, 执行在Qt Creator由项目文件夹构建出的segmentation.exe程序即可。

请注意确保window.cpp中调用loadDict函数时所传路径为您的电脑中dict.txt文件的绝对路径。

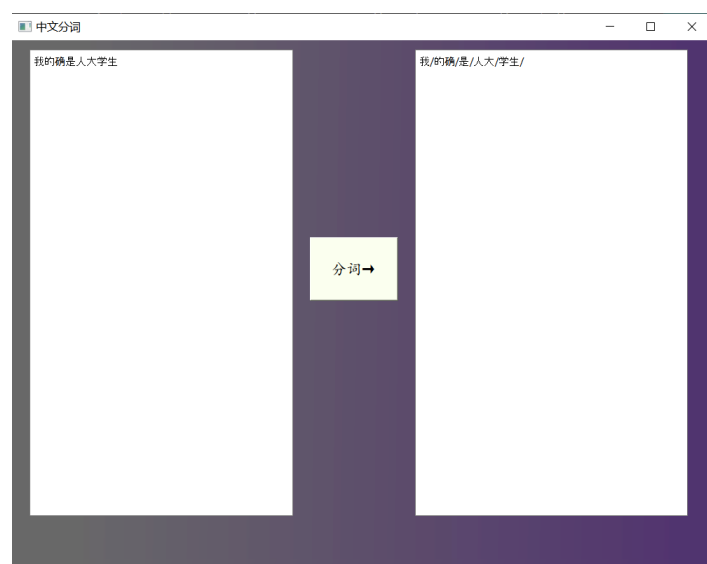
- 程序运行后, 您将看到如下交互界面:



- 界面左侧设为输入框, 右侧设为输出框, 在输入框输入您待分词的由中文字符构成的文本 (长度上限约一万字), 单击界面中间的**分词**按钮, 即可在右侧的输出框看到分词结果。

6. 测试结果

(1) 单句分词测试



(2) 长文本分词测试

中文分词

文建策: 奥陶无人区
尊敬的潘云鹤院士、怀进鹏书记、杜占元局长、李萌副部长、高润生副司长、黄铁军院长、刘伟校长、张磊校友、

各位嘉宾、老师们、同学们:

春水初生, 春林日盛, 春晖方暖。今天我们相聚中国人民大学, 共同见证高瓴人工智能学院的成立。感谢学校和社会各界在学院筹备过程中的大力支持!

习近平总书记在中南海政治局第九次集体学习中指出, 要“奥陶人工智能科技前沿的无人区, 努力在人工智能发展方向和理论、方法、工具、系统等方面取得突破性、颠覆性突破”。今年2月份春节后我们向潘云鹤院士汇报高瓴人工智能学院筹备工作的时候, 他也向我们提出了学院建设要“奥陶无人区”的愿望。

我们的愿景是打造一个能够影响和塑造未来人工智能时代的世界一流学院, 为中国和全球思考并创造“智能而有温度”的未来。在人民大学这样一所以人文社会科学见长的著名学府中建设一个世界一流的人工智能学院是一个全新的挑战。我们需要奥陶的无人区。

首先, 我们要在学院体制机制建设上开拓创新。在学院筹备之初, 学校领导和张磊校友就提出了“高起点、高水平、国际性、创新性”的发展方针, 希望我们大胆探索学院学科建设的一流创新体制机制。高瓴人工智

分词→

文/绝/康/奥陶/无人区/尊敬/的/潘/云/鹤/院士/怀/进/鹏/书记/杜/占/元/局长/李/萌/副/部/长/高/润/生/副/司/长/黄/铁/军/院/长/刘/伟/校/长/张/磊/校/友/各位/嘉/宾/老/师/们/同/学/们/春/水/初/生/春/林/日/盛/春/晖/方/暖/今/天/我/们/相/聚/人/民/大/学/共/同/见/证/高/瓴/人/工/智/能/学/院/的/成/立/感/谢/学/校/和/社/会/各/界/在/学/院/筹/备/过/程/中/的/大/力/支/持/习/近/平/总/书/记/在/中/央/政/治/局/第/九/次/集/体/学/习/中/指/出/要/奥/陶/人/工/智/能/科/技/前/沿/的/无/人/区/努/力/在/人/工/智/能/发/展/方/向/和/理/论/方/法/工/具/系/统/等/方/面/取/得/突/破/性/颠/覆/性/突/破/今/年/月/份/春/节/后/我/们/向/潘/云/鹤/院/士/汇/报/高/瓴/人/工/智/能/学/院/筹/备/工/作/的/时/候/他/也/向/我/们/提/出/了/学/院/建/设/要/奥/陶/无/人/区/的/期/望/我/们/的/愿/景/是/打/造/一/个/能/够/影/响/和/塑/造/未/来/人/工/智/能/时/代/的/世/界/一/流/学/院/为/中/国/和/全/球/思/考/并/创/造/智/能/而/有/温/度/的/未/来/在/人/民/大/学/这/样/一/所/以/人/文/社/会/科/学/见/长/的/著/名/学/府/中/建/设/一/个/世/界/一/流/的/人/工/智/能/学/院/是/一/个/全/新/的/挑/战/是/我/们/需/要/奥/陶/的/无/人/区/首/先/我/们/要/在/学/院/体/制/机/制/建/设/上/开/拓/创/新/在/学/院/筹/备/之/初/学/校/领/导/和/张/磊/校/友/就/提/出/了/高/起/点/高/水/平/国/际/性/创/新/性/的/发/展/方/针/希/望/我/们/大/胆/探/索/学/院/学/科/建/设/的/一/流/创/新/体/制/机/制/高/瓴/人/工/智/能/学/院/将/积/极/探/索/新/的/科/研/和/人/才/评/价/制/度/不/拘/一/格/选/拔/和/培/养/青/年/人/才/学/院/将/借/鉴/国/际/名/校/规/范/创/新/人/才/管/理/机/制/教/师/序/列/待/全/面/实/行/制/度/学/院/将/为/教/师/打/造/国/际/一/流/的/工/作/生/活/环/境/提/供/与/国/际/知/名/大/学/和/研/究/机/构/接/轨/的/学/术/平/台/服/务/以/平/作/育/多/位/优/秀/人/才

中文分词

吴志成: 从微软到人大
从企业转到高校, 就像过去人生中的许多重要转折点一样, 吴志成并不是刻意追求, 而是顺其自然。但他对凡事均竭尽全力, 总能在看似被动的经历中取得佳绩: 在过去的六年里, 他先是担任人大信息学院特别研究员, 随后被破格晋升为教授, 如今又被任命为人大高瓴人工智能学院副院长。

事实上, 作为一路拔尖的学霸, 从小学到初中、高中, 甚至是大学, 吴志成的学习成绩一直名列前茅, 稳扎稳打。但高考后, 吴志成没有选择报考湖北, 而是去了南开读计算机。

当时, 学校负责为学生提供报考意见的老师跟他谈, 前一年的承德市状元去了南开大学读计算机, 体验非常好。

那时候, 计算机刚刚起步, 周边人一致认为这个专业很有前途。作为理科生, 吴志成的内心对前沿技术自然十分向往, 于是便报了南开大学信息技术科学学院计算机专业。

结果, 吴志成去了南开后发现: “那个师兄报的不是计算机, 而是国贸专业。”真是误打误撞。

更让吴志成没有想到的是, 因为一个“信息不对称”, 他将在南开度过最宝贵的九年时光: 本科、保研、转博, 成为一名中国“土生土长”的计算机博士, 重点研究领域还是信息检索。

分词→

吴/志/成/从/微/软/到/人/大/从/企/业/转/到/高/校/就/像/过/去/人/生/中/的/许/多/重/要/转/折/点/一/样/吴/志/成/并/不/是/刻/意/追/求/而/是/顺/其/自/然/但/他/对/凡/事/均/竭/尽/全/力/总/能/在/看/似/被/动/的/经/历/中/取/得/佳/绩/在/过/去/的/六/年/里/他/先/是/担/任/人/大/信/息/学/院/特/别/研/究/员/随/后/破/格/晋/升/为/教/授/如/今/又/被/任/命/为/人/大/高/瓴/人/工/智/能/学/院/副/院/长/事/实/上/作/为/一/路/拔/尖/的/学/生/成/绩/一/直/名/列/前/茅/稳/扎/稳/打/但/高/考/后/吴/志/成/没/有/选/择/报/考/湖/北/而/是/去/了/南/开/读/计/算/机/当/时/学/校/里/负/责/为/学/生/提/供/报/考/意/见/的/老/师/跟/他/谈/前/一/年/的/承/德/市/状/元/去/了/南/开/大/学/读/计/算/机/体/验/非/常/好/那/时/候/计/算/机/刚/刚/起/步/周/边/人/一/致/认/为/这/个/专/业/很/有/前/途/作/为/理/科/生/吴/志/成/的/内/心/对/前/沿/技/术/自/然/十/分/向/往/于/是/便/报/了/南/开/大/学/信/息/技/术/科/学/学/院/计/算/机/专/业/结/果/吴/志/成/去/了/南/开/发/现/那/个/师/兄/报/的/不/是/计/算/机/而/是/国/贸/专/业/真/是/误/打/误/撞/更/让/吴/志/成/没/有/想/到/的/是/因/为/一/个/信/息/不/对/称/他/将/在/南/开/度/过/最/宝/贵/的/九/年/时/光: 本/科/保/研/转/博/成/为/一/名/中/国/土/生/土/长/的/计/算/机/博/士/重/点/研/究/领/域/还/是/信/息/检/索/

中文分词

张磊: 选择做时间的朋友
希望大家都能坚持自己内心的选择, 不骄不躁, 好故事都是来自于有挑战的生活; 持之以恒, 时间终将会成为你的朋友。在这里与同学们共勉。

除了选择的几个问题。作为人大的校友, 我觉得还有一点很重要, 那就是我们不仅要掌握科学思维的能力, 还要心中长存人文精神的火种。当今时代, 伴随基因技术、机器人和人工智能技术的发展, 科技爆炸、奇点临近, 人类将进入新纪元。我们的生活也会迎来巨大的挑战。而大家在人大的学习生活, 恰恰赋予了我们广阔的视野和人文情怀。这得帮助我们处乱不惊、面对纷繁复杂的世界, 不断去追问问题的本质。我本科是学国际金融的, 没有编过也没有技术背景。但是我后来照样投资了一批科技企业, 它们现在在各自领域内引领世界潮流。我感谢咱们人大的人文教育, 相信同学们也会从中获益无穷。

作为投资人, 我常说起我的三个投资哲学: “守正用奇”、“弱水三千, 但取一瓢”、“桃李不言, 下自成蹊”, 这三句话分别源自《道德经》、《论语》、《史记》。虽然现代金融投资的工具和方法大多源于西方, 但如何使用好这些工具, 我还是要感谢我们优秀的中国哲学思想和传统民族文化。我们要有充分的文化自信, 要珍惜人大给我们的“人文土壤”, 好好汲取营养。过去未去, 未来已来。我希望我们人大学子, 以后不管进入哪个行业/从事什么工作/都能保持乐观和激情, 用人文的情怀去雕琢自己/美化身边/赠人玫瑰, 手有余香。

分词→

张/磊/选/择/做/时/间/的/朋/友/希/望/大/家/都/能/坚/持/自/己/内/心/的/选/择/不/骄/不/躁/好/故/事/都/是/来/自/于/有/挑/战/的/生/活/持/之/以/恒/时/间/终/将/会/成/为/你/的/朋/友/在/这/里/与/同/学/们/共/勉/除/了/选/择/的/几/个/问/题/作/为/人/大/的/校/友/我/觉/得/还/有/一/点/很/重/要/那/就/是/我/们/不/仅/要/掌/握/科/学/思/维/的/能/力/还/要/心/中/长/存/人/文/精/神/的/火/种/当/今/时/代/伴/随/基/因/技/术/机/器/人/和/人/工/智/能/技/术/的/发/展/科/技/爆/炸/奇/点/临/近/人/类/将/进/入/新/纪/元/我/们/的/生/活/也/会/迎/来/巨/大/的/挑/战/而/大/家/在/人/大/的/学/习/生/活/恰/恰/赋/予/了/我/们/广/博/的/视/野/和/人/文/情/怀/这/得/帮/助/我/们/处/乱/不/惊/面/对/纷/繁/复/杂/的/世/界/不/断/去/追/问/问/题/的/本/质/我/本/科/是/学/国/际/金/融/的/没/有/编/过/程/序/也/没/有/技/术/背/景/但/是/我/后/来/照/样/投/资/了/一/批/科/技/企/业/它/们/现/在/在/各/自/领/域/内/引/领/世/界/潮/流/我/感/谢/咱/们/人/大/的/人/文/教/育/相/信/同/学/们/也/会/从/中/获/益/无/穷/作/为/投/资/人/我/常/说/起/我/的/三/个/投/资/哲/学/守/正/用/奇/弱/水/三/千/但/取/一/瓢/桃/李/不/言/下/自/成/蹊/这/三/句/话/分/别/源/自/道/德/经/论/语/史/记/虽/然/现/代/金/融/投/资/的/工/具/和/方/法/大/多/源/于/西/方/但/如/何/使/用/好/这/些/工/具/我/还/是/要/借/鉴/我/们/优/秀/的/中/国/哲/学/思/想/和/传/统/民/族/文/化/我/们/要/珍/惜/人/大/给/我/们/的/人/文/土/壤/好/好/汲/取/营/养/过/去/未/去/未/来/已/来/我/希/望/我/们/人/大/学/子/以/后/不/管/进/入/哪/个/行/业/从/事/什/么/工/作/都/能/保/持/乐/观/和/激/情/用/人/文/的/情/怀/去/雕/琢/自/己/美/化/身/边/赠/人/玫/瑰/手/有/余/香/