



大数据智能之道法术

Big Data Intelligence: Tao, Principle and Tactics

张华平 博士



大数据搜索与挖掘实验室

kevinzhang@bit.edu.cn

www.nlpir.org

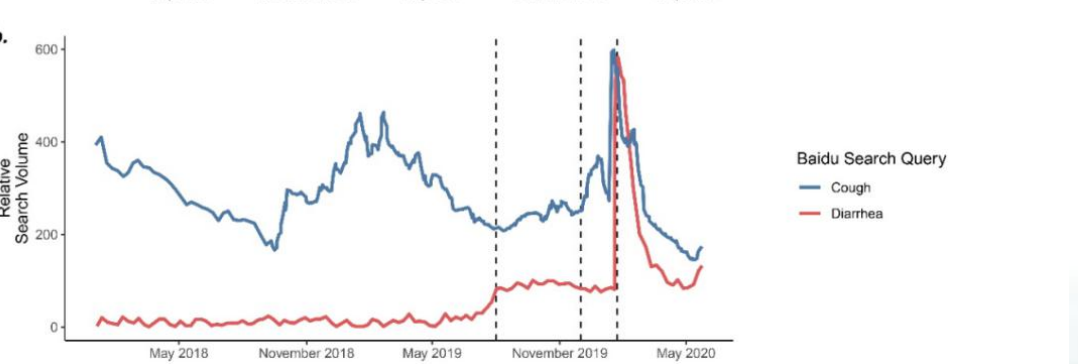
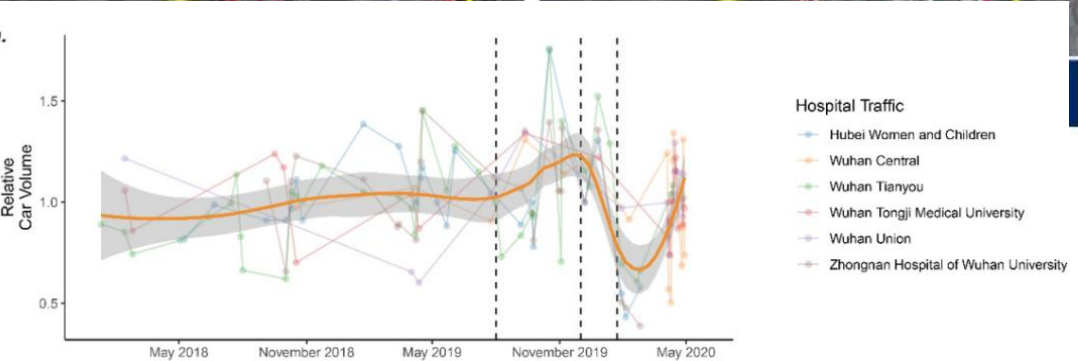
2021.9



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

大数据智能：新冠预测与五角大楼停车场指数

NEWS EXCLUSIVE HUBEI WOMEN AND CHILDREN HOSPITAL



2012年10月13日这一天，西南侧车辆突然增多约100辆左右

10月13日美海军"蒙彼利埃"号潜艇与提康德罗加级巡洋舰"圣哈辛托"号在东部海域相撞



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

从棱镜手机监控看大数据...



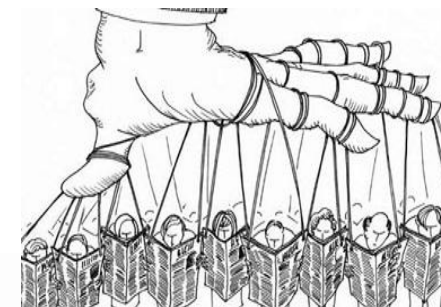
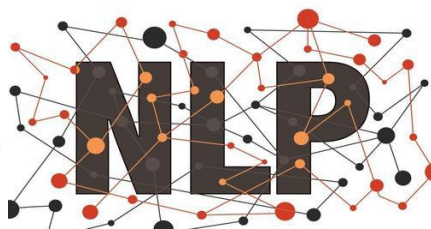
工 大 学

BEIJING INSTITUTE OF TECHNOLOGY

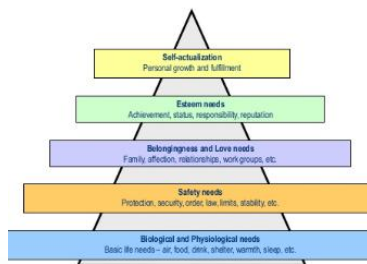
社交媒体大数据背景下的信息爆炸与社群操控



NATURAL LANGUAGE PROCESSING



Maslow's Hierarchy of Needs Model



Cambridge Analytica



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

国家安全问题:分裂暴恐势力藏身社交媒体大数据

港独台独勾结，危害国家稳定。境外势力插手，妄图独立台湾，



社会治理问题：违法犯罪信息变身逃避常规扫描

您好：
 常年为各大公司企业代开增值税专用以及增值税普通 发票 ，可对公走账签订合同，工程款，材料款，广告费，制作费，发布费，服务费，劳务费，咨询费，等等各行业正轨发票，验证后付款。
 卢经理：13671133317
 或添加微信号：fapiao13671133317
 希望能与朋友们建立长久的合作关系，此信息长期有效有需要的朋友们请保留备用

嫌去澳门太麻烦，在家开户就能玩
 斥资1500亿打造网上最佳品牌，持
 菲律宾正规博彩执照，单用户日最
 高可提款5000万，3-5分钟到账！

✅【帝王国际娱乐城】
 官方网址：111v.cc
 📦 开户首存1000送500
 一直被模仿，从未被超越！
 专营：百家乐·电子游戏·体育赛事

专注互联网彩票公司
 ✅（速彩购彩中心）
 官方网址：999y1.com
 📦 收到短信3天内开户即送188彩金
 主营：各种互联网彩票，真人视讯

你好 ftaqj
 为了您的合理避.税，减轻企业负担，为您提供 D.K 各类『发*票』业务。

电话：135*3412*6969 另加微VX+ A-33680

负责人，您好
 我司有各地公司剩余票/据凭/证，正/规有效，主要有各类销售行业，材料费,工/程费，安装，办公费，技术服/务费,住/宿费/会务费，设计费，咨询/培训费等各种费用，点低，诚信！
 联//系人：成先生
 手//机：136-1288-5985
 在线加Q：502603743

01

假发票短信泛滥

02

涉赌短信屡禁不绝

03

假发票邮件不断“变种”



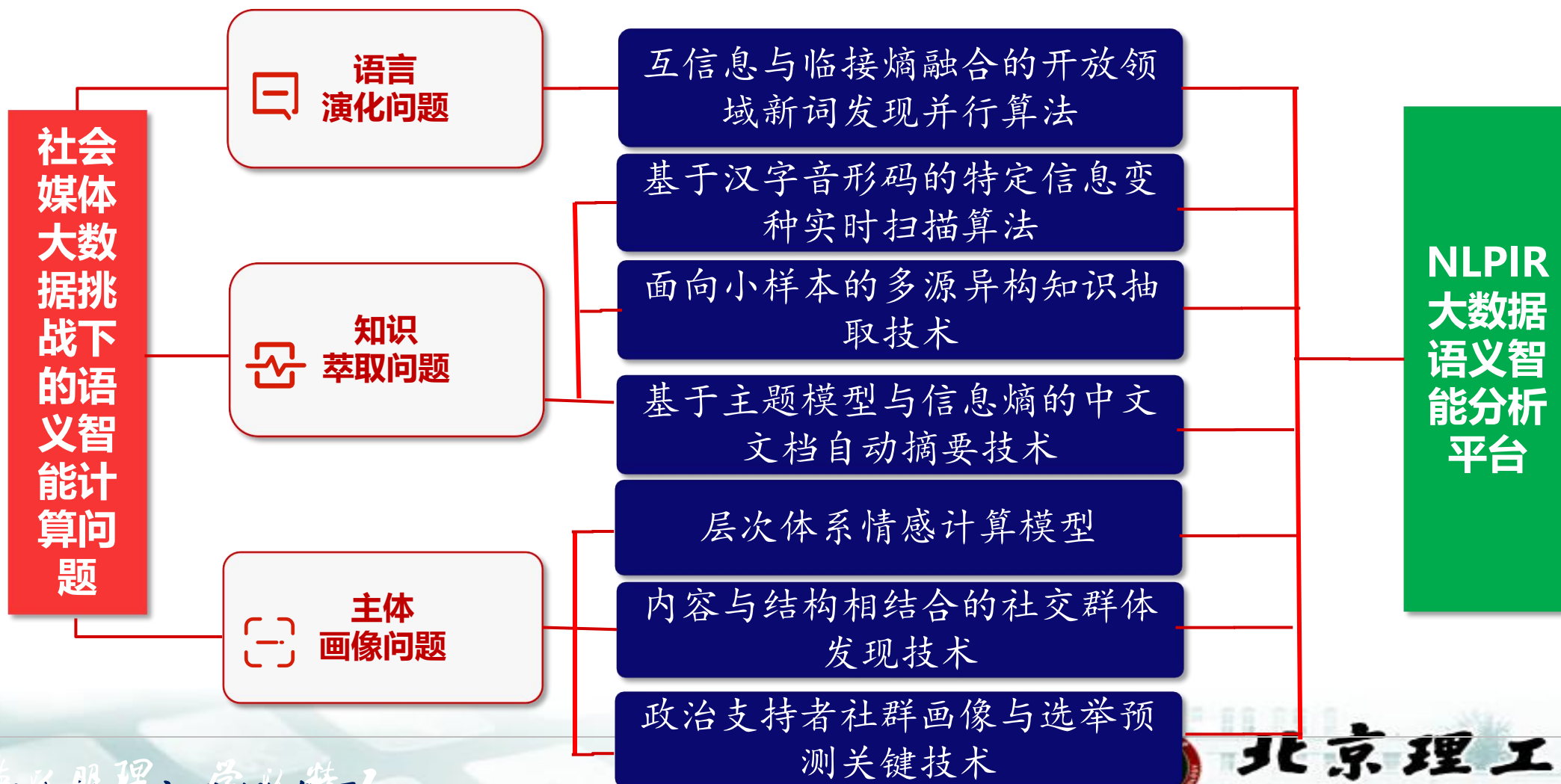
大数据智能分析架构

一个科学问题

三大技术挑战

七项创新技术

一个集成平台



大数据 人工智能 自然语言处理

推动互联网、大数据、人工智能
和实体经济深度融合

-十九大报告

人工智能，现代科学皇冠上的明珠；
自然语言处理则是人工智能皇冠上的
明珠。

-周明，ACL主席

人工智能的突破在自然语言理解，懂
语言者得天下

-沈向洋，原微软全球执行副总裁

数据

信息

知识

情报

认知智能

- 认知能力：与人的交流、交互与交融
- 自然语言处理、知识推理

感知智能

- 感知能力：受限的环境
- 自动驾驶、人脸识别、传感器等

计算智能

- 计算能力：人工定义的严格规则
- AI剪枝优化决策，大数据存储与计算

什么是大数据智能

➤ 我们的见解：

- 大数据智能是指从客观存在的全量超大规模、多源异构、实时变化的微观数据中，利用自然语言处理、信息检索、机器学习等技术抽取知识，转化为智慧的方法学。
- 神即道，道即法，道法自然，如来
- 智能为道、数据为法、语义为术



大数据智能将带来巨大的变革与机会

AI'S TRANSFORMATIVE POTENTIAL

Projected Global Economic Effects of AI by 2030



CHINA'S AI ASPIRATIONS

China aims to lead the world in AI technologies by 2030.

The Chinese government aims to build a US\$15 billion AI market by 2018.

However, by 2030, AI will provide an expected 26% boost to GDP.

(Source: CNBC, Technode, PwC)

15.7 万亿美元
2030年

PWC预测的世界AI市场规模
(2017年GDP 中国+印度=15.3万亿美元)





认识大数据



什么是大数据

- ➔ **Wiki:** **Big data** is the term for a collection of data sets so large and complex that it becomes difficult to process using on-hand database management tools or traditional data processing applications.
- ➔ 维克托 《大数据时代》：大数据指不用随机分析法（抽样调查）这样的捷径，而采用所有数据的方法。



什么是大数据智能

➔ 我们的见解:

- 大数据智能是指从客观存在的全量超大规模、多源异构、实时变化的微观数据中，利用自然语言处理、信息检索、机器学习等技术抽取知识，转化为智慧的方法学。
- 是一场新的科技革命，也是思想方法的革命。（全量分析，让数据说话；承认并客观地认识世界的混杂性；相关性挖掘替代因果推断）





这是显微镜下的世界



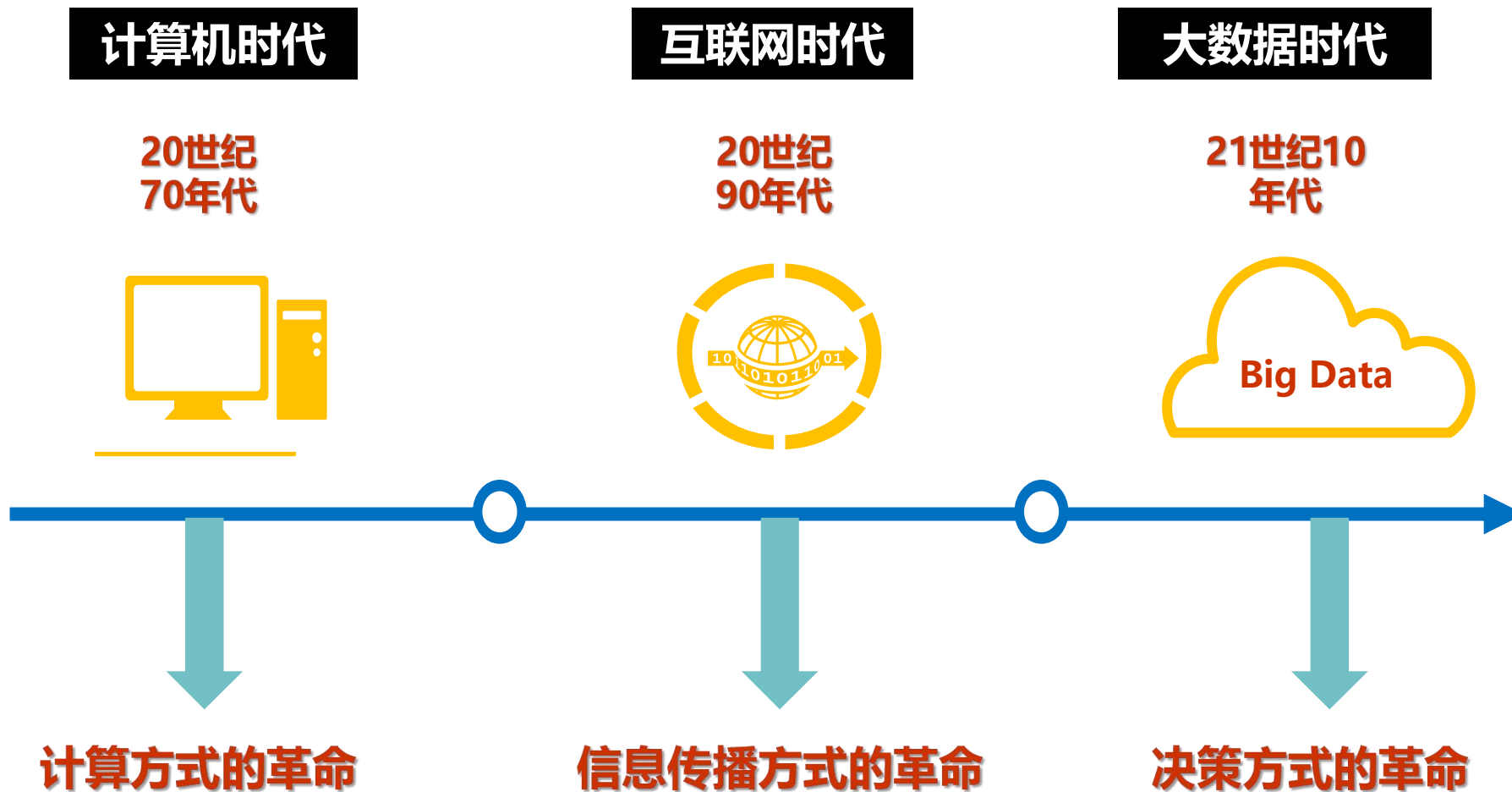
这是望远镜中的宇宙



杨达才启示：1+1>>2才是大数据



近半世纪来的三次革命



大数据颠覆决策模式

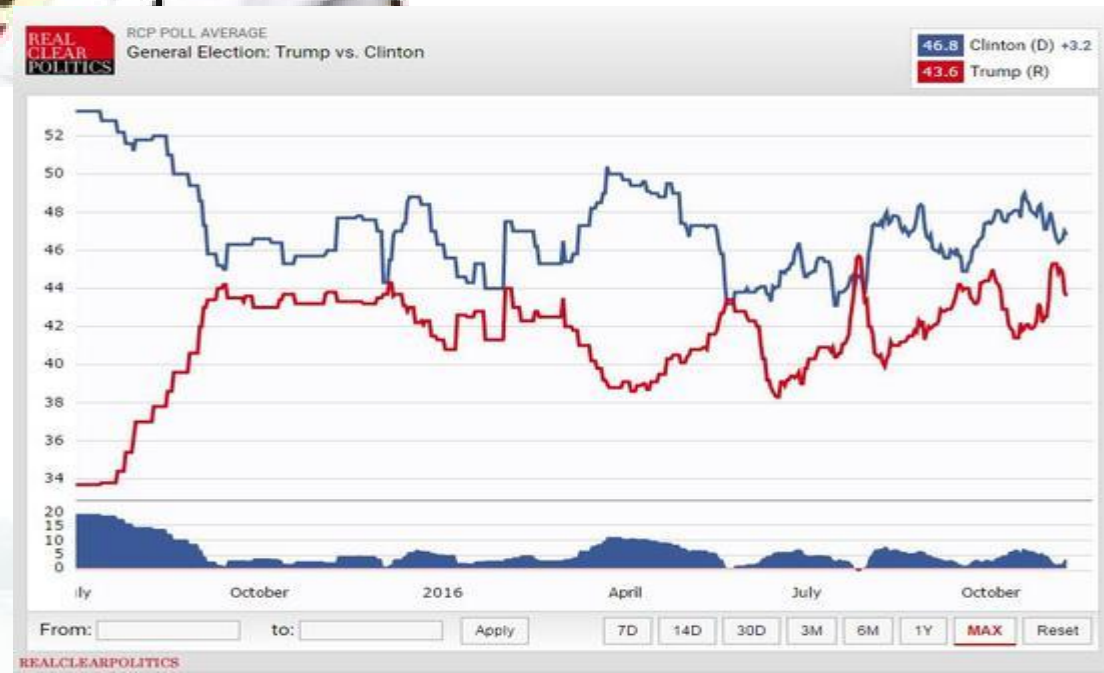


ation

Ultimate

北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

小数据精英 VS 大数据庶民

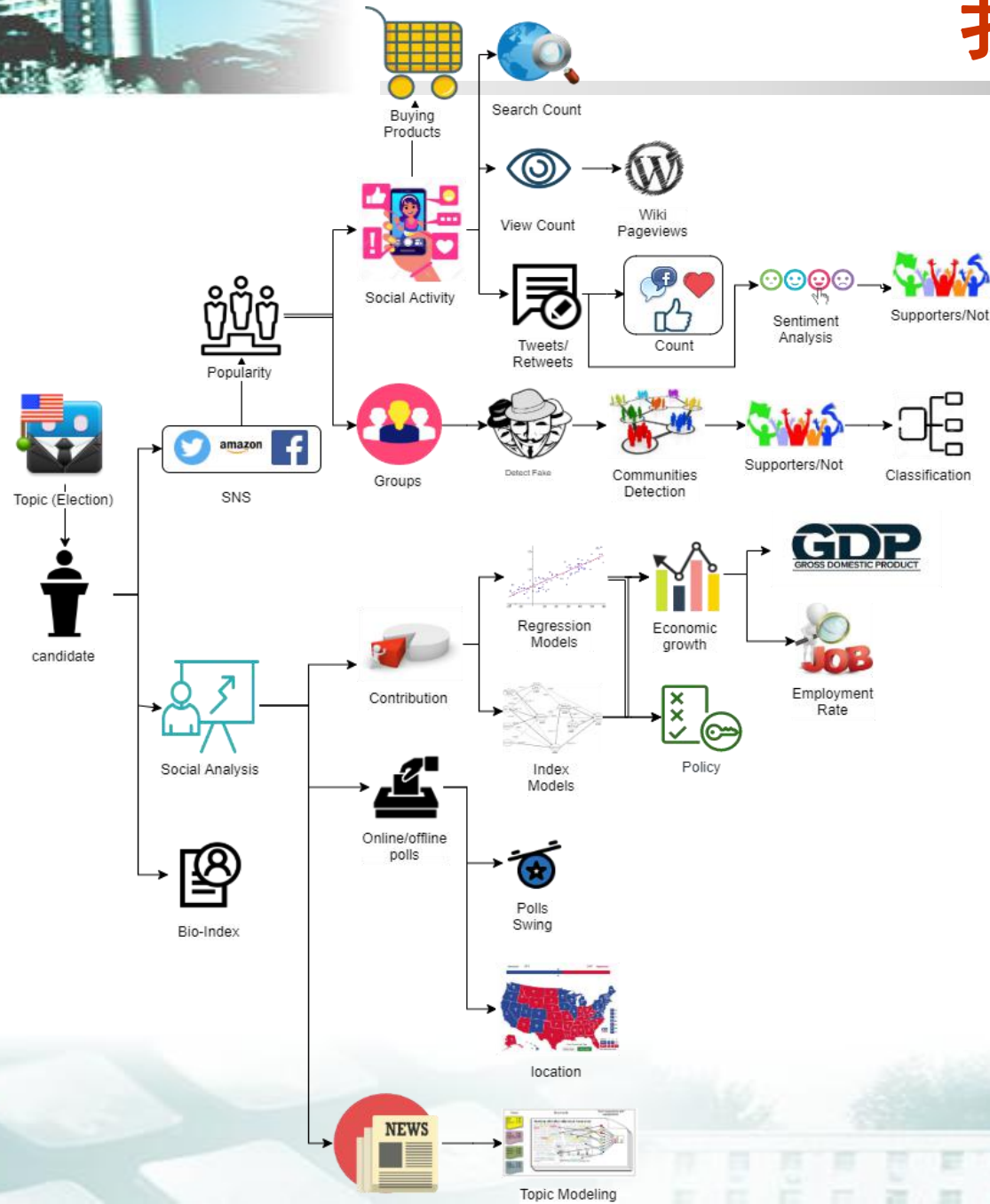




小数据精英 VS 大数据庶民



我们的大数据预测模型



➤ 2018年提前半年成功预测巴基斯坦大选；

➤ 2020年初成功预测蔡英文连任，预测的得票率误差在5%；

➤ 2020年预测美国大选特朗普胜算大于拜登10%，目前失败但仍然颠覆了民调结果。



幸存者偏差：你不了解的中国

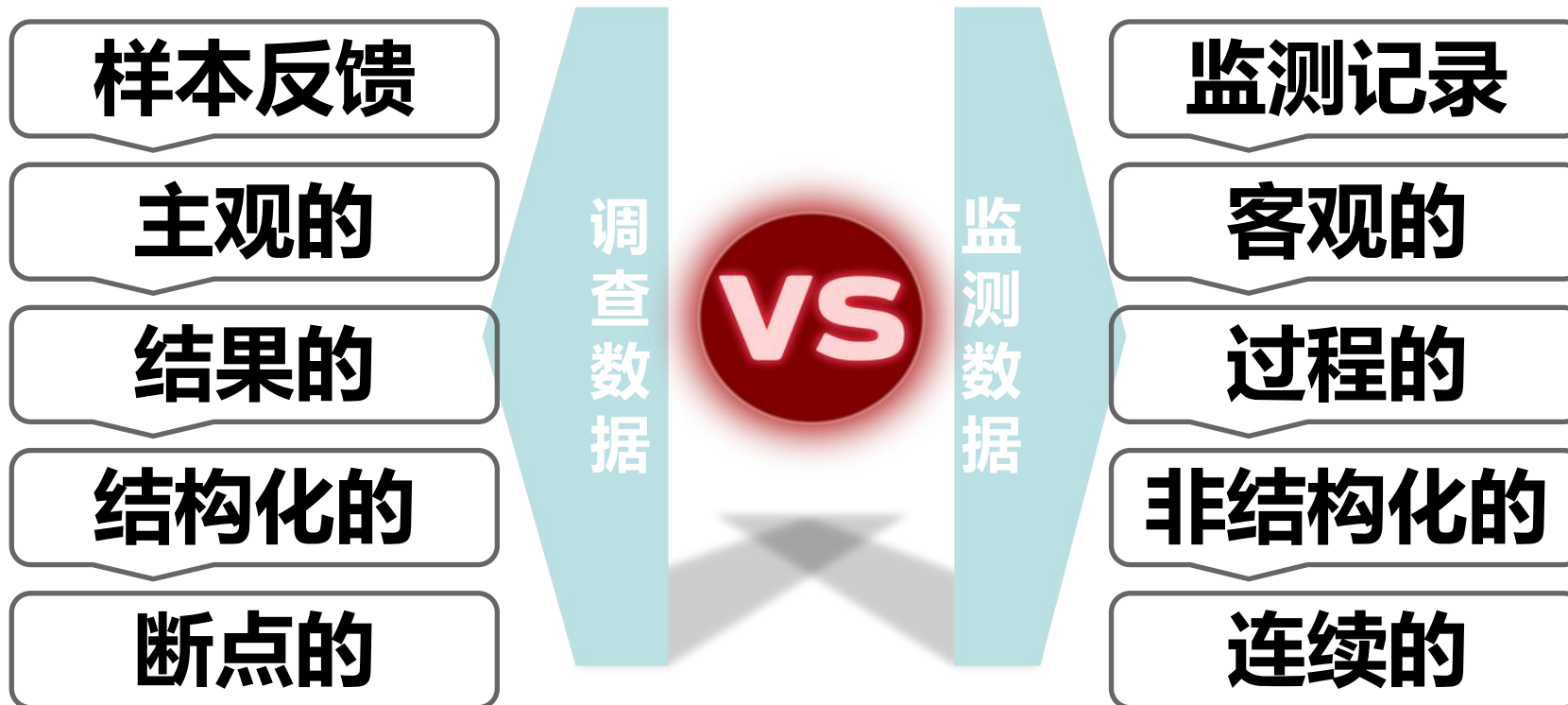
- 10亿人没有坐过飞机；4亿人没有使用过抽水马桶
- 6亿人次出国，出国的人数只占5%
- 2019年6.18，消费升级，市场下沉，四五线贡献率超过7成
- 618奇瑞小蚂蚁400 1分钟销量387台，比亚迪7805台
- 7.14高炮背后的庞大群体



大数据时代的特征



小数据和大数据的区别





WHAT IS A.I.?

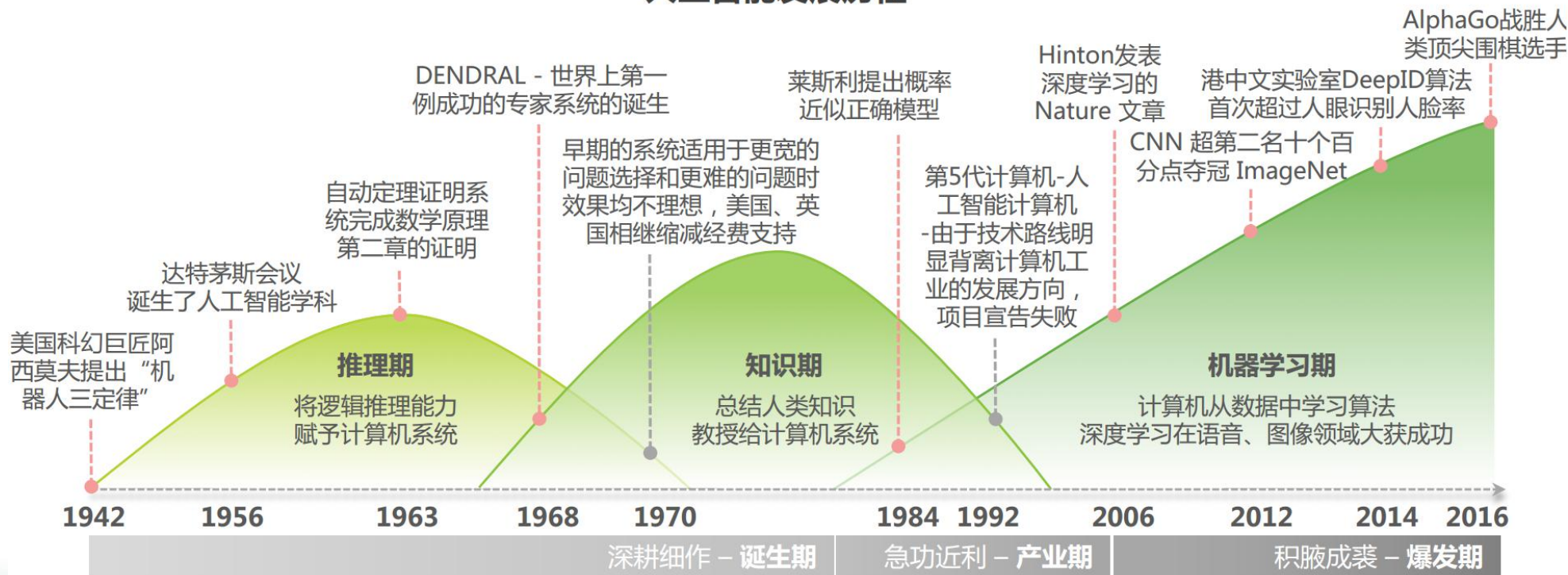
认识人工智能



人工智能定义与发展历程

➔ 1956年，约翰麦卡锡(John McCarthy)在达特茅斯会议上首次提出人工智能（Artificial Intelligence: AI）的定义：使一部机器的反应方式像一个人在行动时所依据的智能。

人工智能发展历程

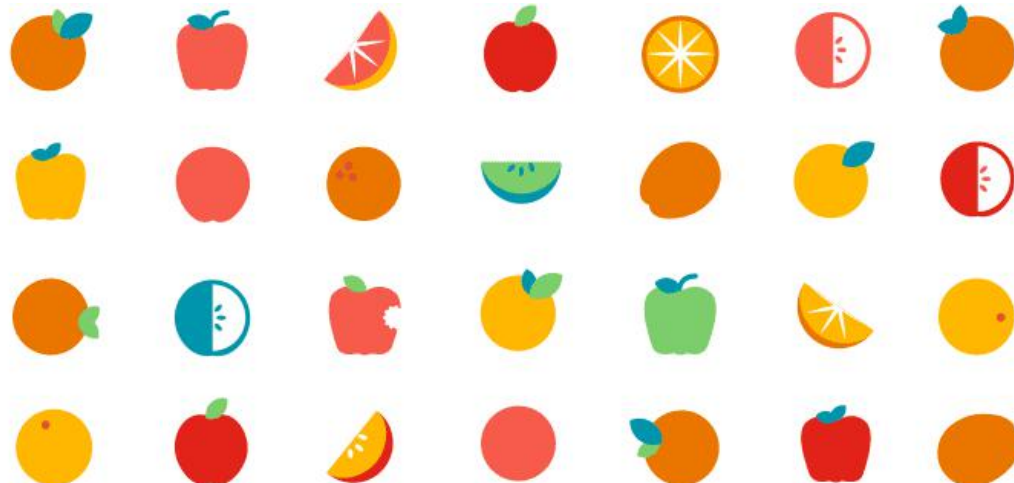
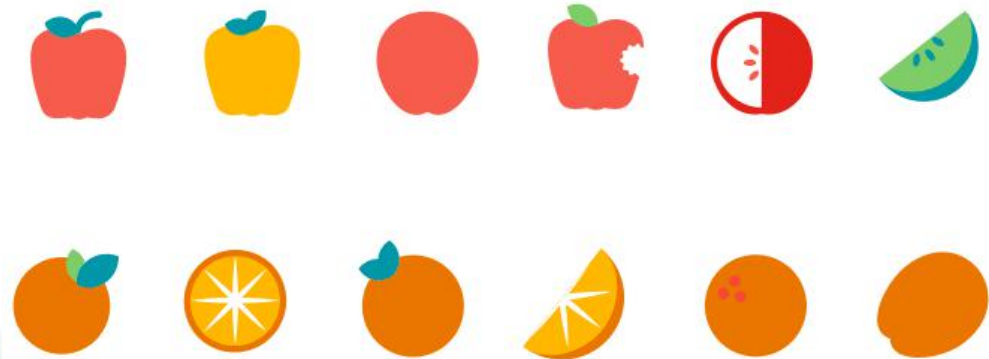
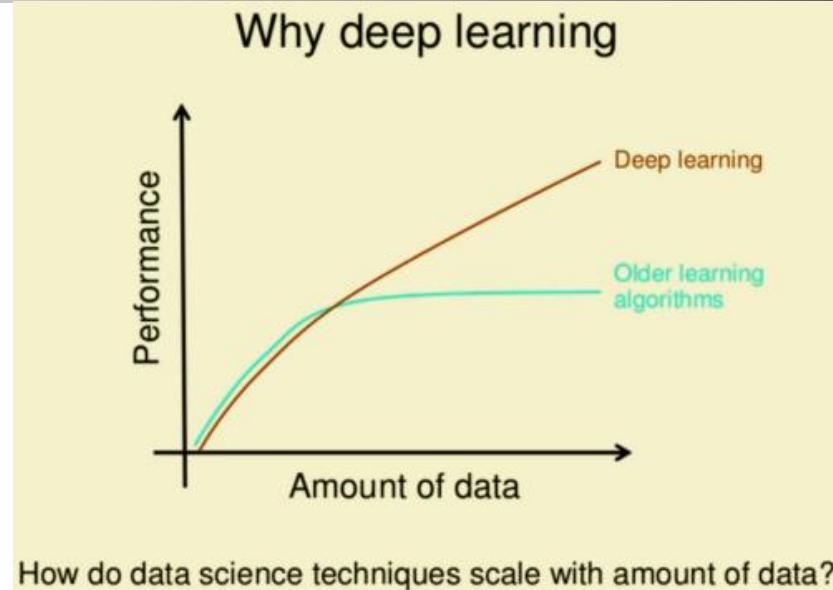


注：图来自于艾瑞收集整理



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

机器学习与深度学习



人工智能的三个境界

超人工智能

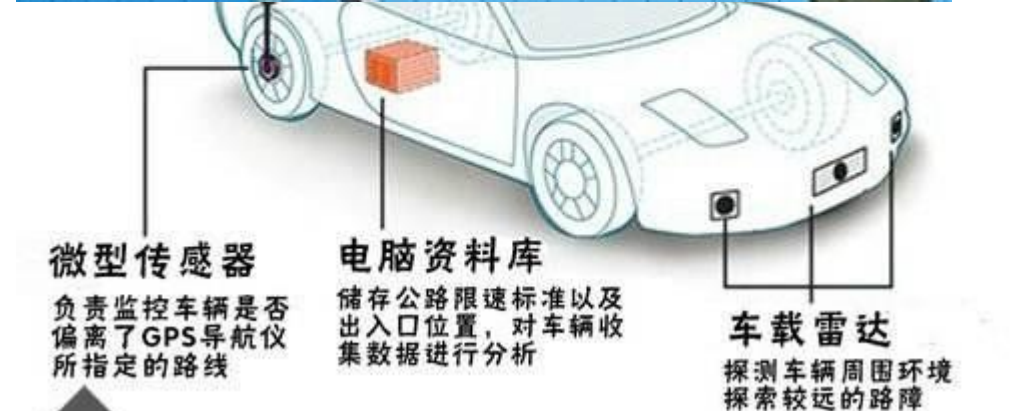
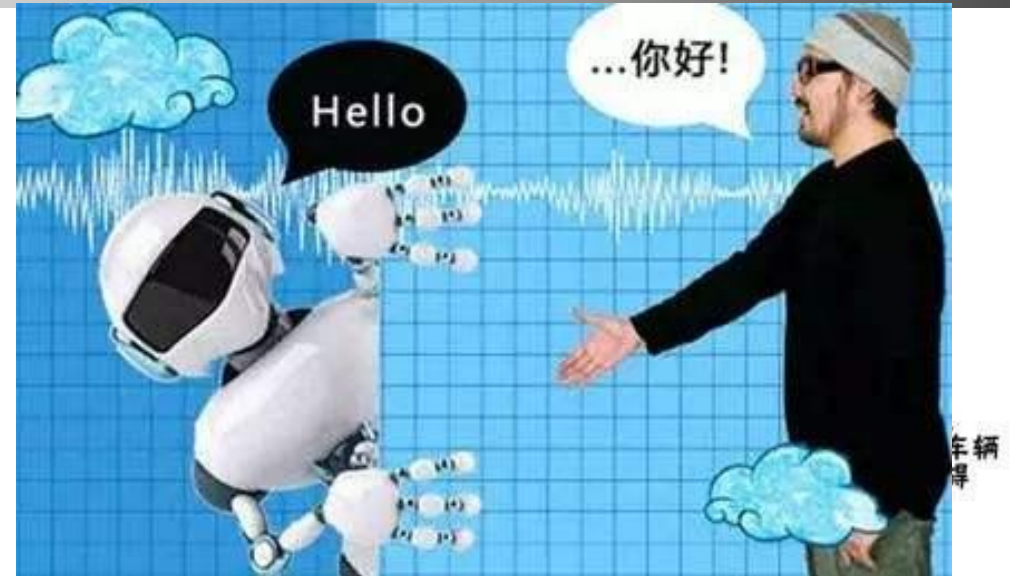
- 认知智能：与人的交流、交互与交融
- 自然语言理解、知识推理

强人工智能

- 感知智能：受限的环境
- 自动驾驶、人脸识别、传感器等

弱人工智能

- 计算智能：人工定义的严格规则
- AI剪枝优化决策，大数据存储与计算

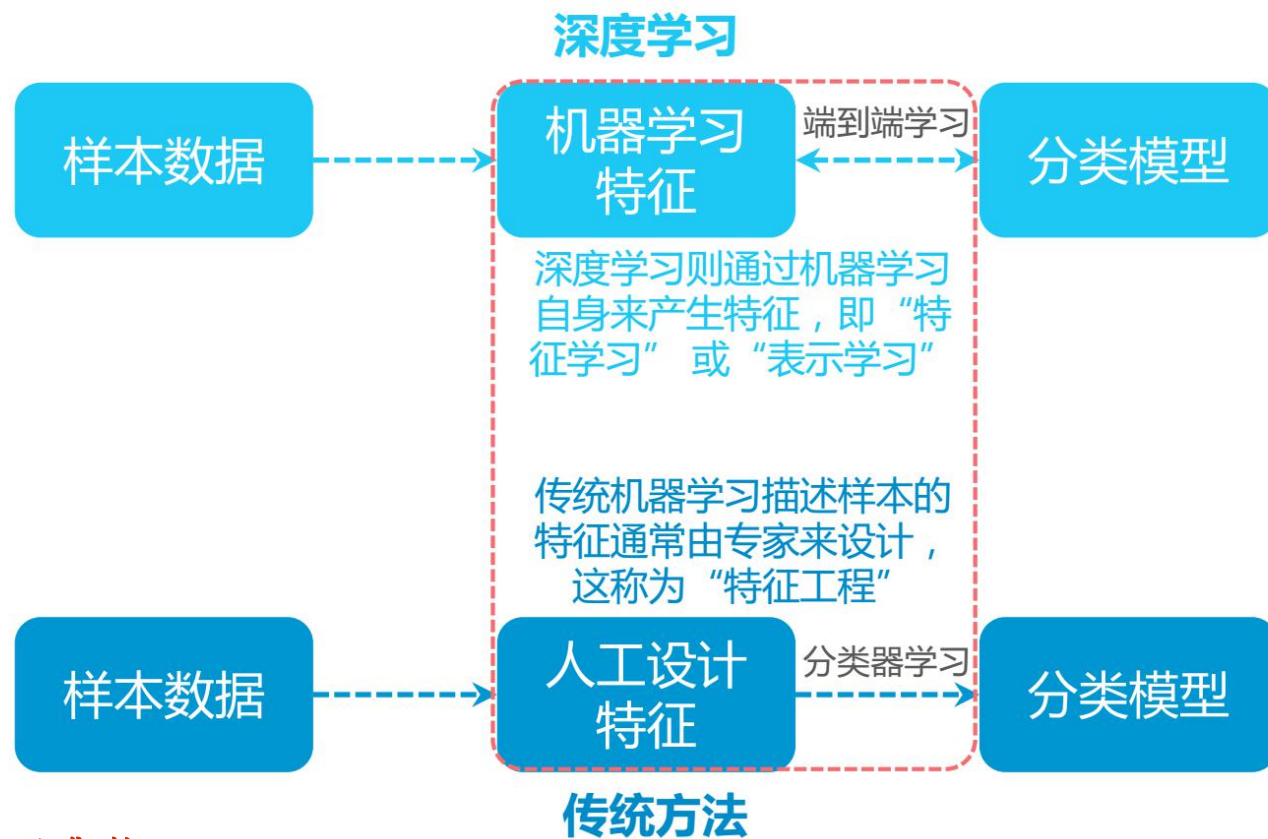


深度学习vs机器学习vs人工智能

深度学习 < 机器学习 < 人工智能



深度学习与传统方法的区别

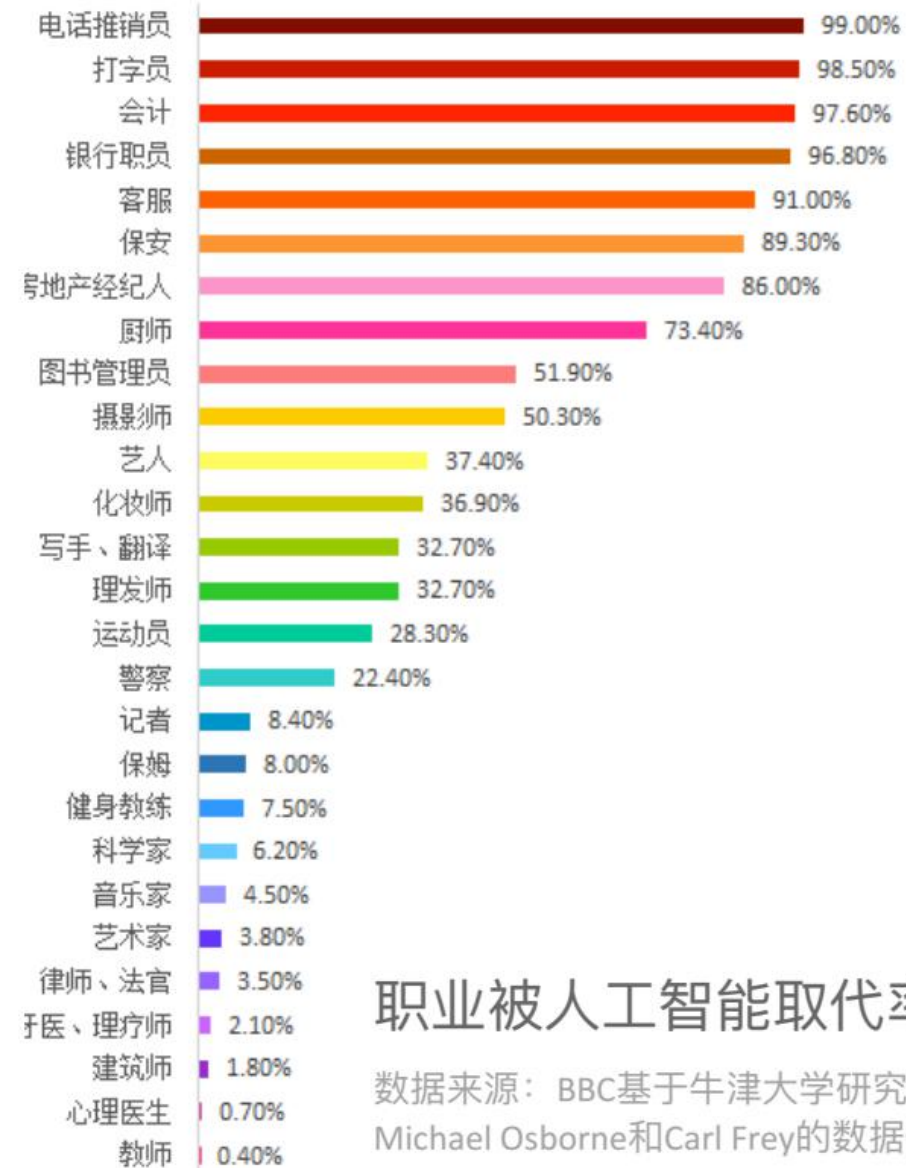


注：图来自于艾瑞收集整理



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

AI将取代大量人力工作



职业被人工智能取代率

数据来源：BBC基于牛津大学研究者
Michael Osborne和Carl Frey的数据体系分析



翻译



记者



司机



客服



保姆



助理



保安



交易员

甚至还有...



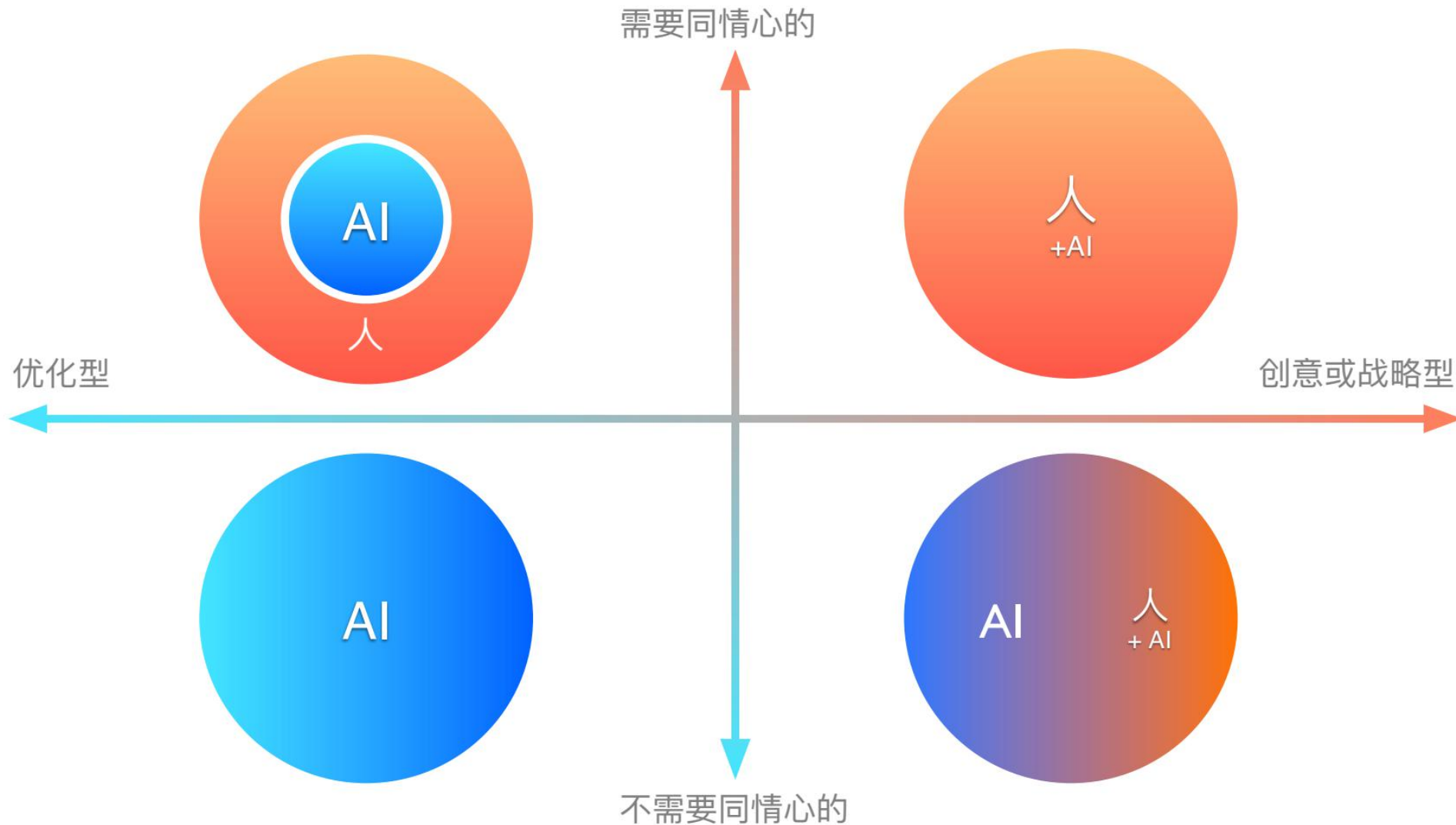
法律咨询师



放射科医生



人工智能替代性分析



人工智能不能干什么？

- 跨领域推理
- 抽象能力
- 审美
- 知其然知其所以然
- 常识推断
- 自我意识
- 情感



巴别通天塔之惑

NLP两大挑战:

歧义消解
人类知识



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

人工智能？

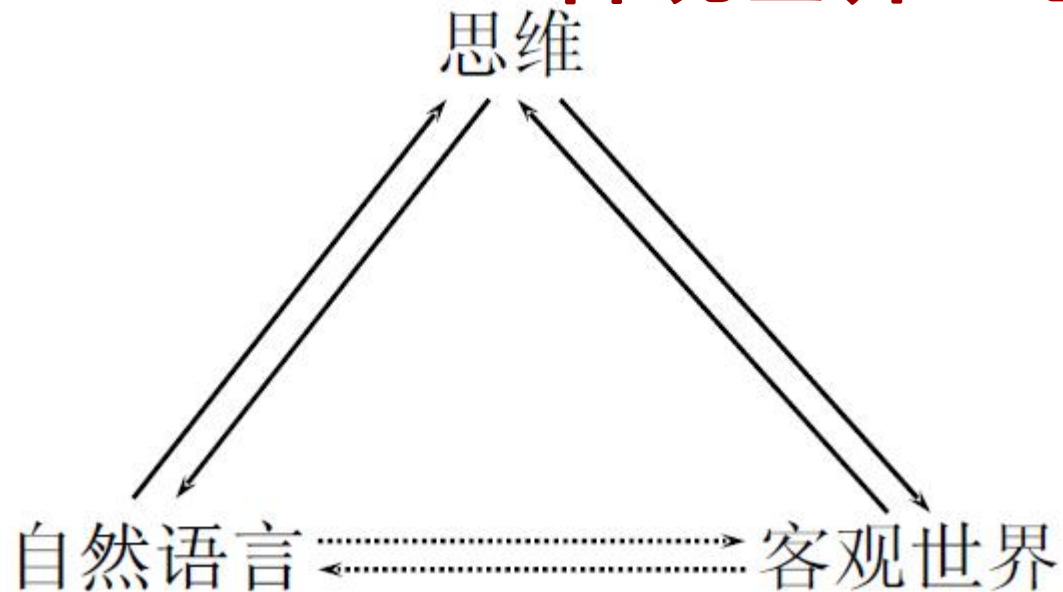
➔ 足球队和乒乓球队：一个谁也打不过，一个谁也打不过。



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY



客观世界->思维->自然语言



➔ 衰减效应：

- 思维最多只能反映80%的客观世界；
- 自然语言只能反映80%的思维：词不达意，答非所问；
- 听众最多只能听懂80%；
- 听懂的部分只有80%能反映到思维中；
- 分析客观世界的最多只能利用80%。



➤ 自然语言处理：技术概念

- is a field of computer science, artificial intelligence and computational linguistics concerned with the interactions between computers and human (natural) languages, and, in particular, concerned with programming computers to fruitfully process large natural language corpora.

➤ 计算语言学：学科概念

- Computational linguistics is an interdisciplinary field concerned with the statistical or rule-based modeling of natural language from a computational perspective.

➤ 文本挖掘：应用概念

- is the process of deriving high-quality information from text.

计算语言学定义

计算语言学是一门以**计算**为手段对**自然语言**进行研究和**处理**的科学。

Computational Linguistics

Natural language processing

"Wherever there is Artificial Intelligence,
there is Artificial Stupidity."

“哪里有人工智能，哪里就有人工愚蠢”。

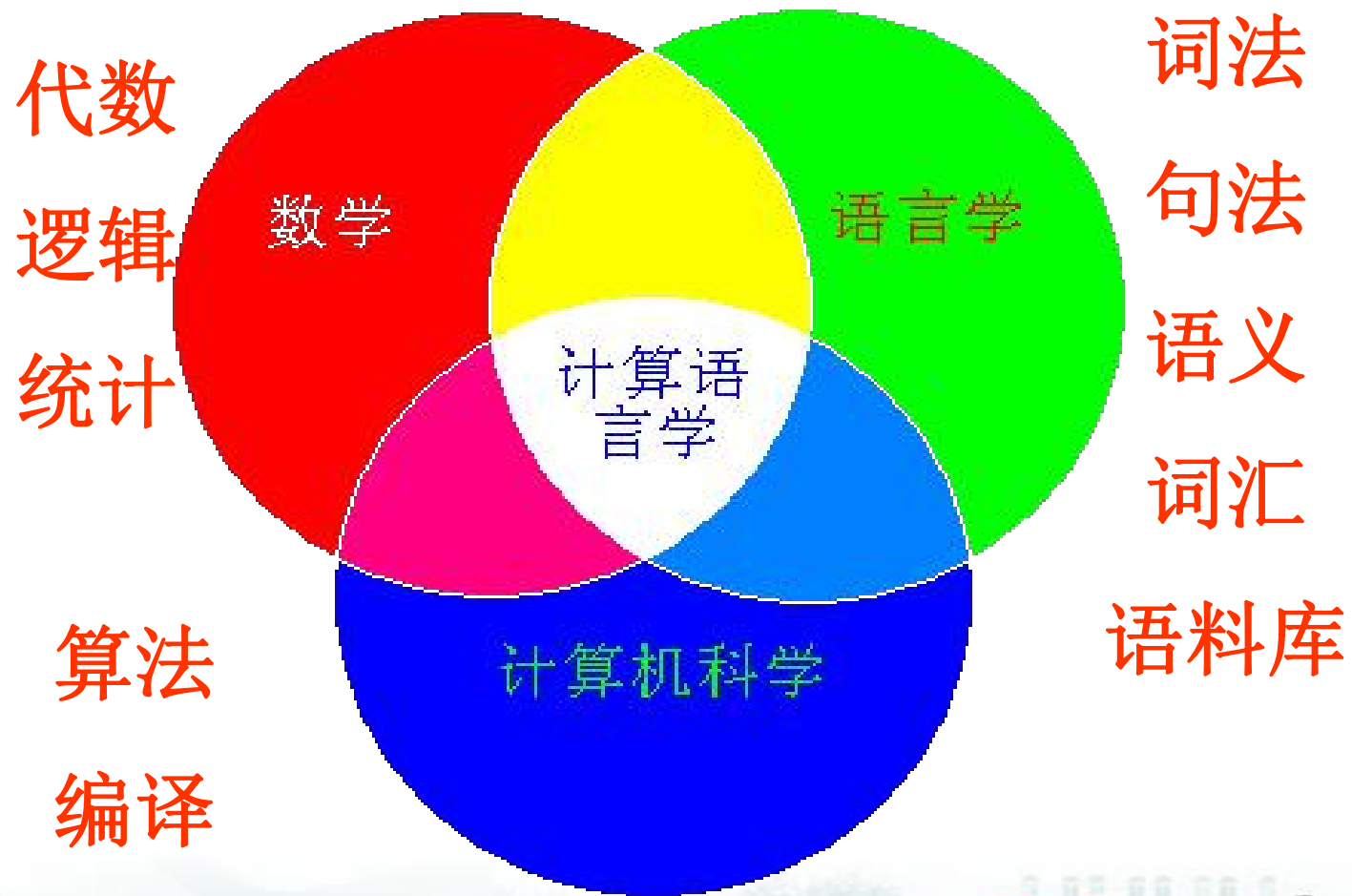
(从一到万, seed→see+ed)



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY



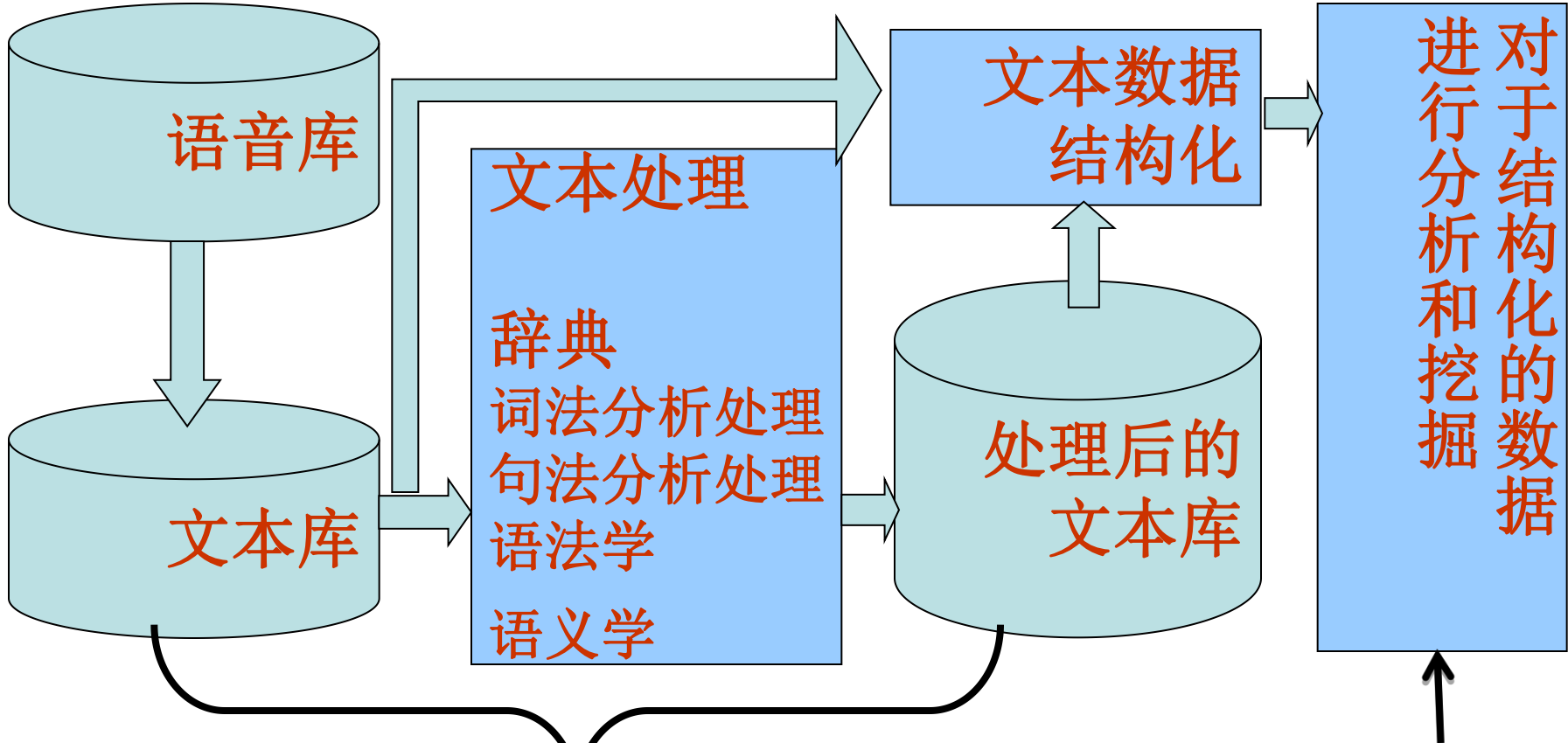
计算语言学与其他学科的关系



人工智能



文本挖掘的框架



自然语言处理

数据挖掘



➤ 基础理论

- 自动机 形式逻辑 统计机器学习, 汉语语言学 形式语法理论

➤ 语言资源

- 语料库 词典

➤ 关键技术

- 汉字编码 词法分析 句法分析 语义分析 文本生成 语音识别

➤ 应用系统

- 人机对话, 机器翻译, 社交网络分析, 分类, 聚类, 检索, 过滤, 信息抽取, 音字转换



NLP主要算法体系-编码识别为例

➤ 理性主义（规则方法）：Rule-Based

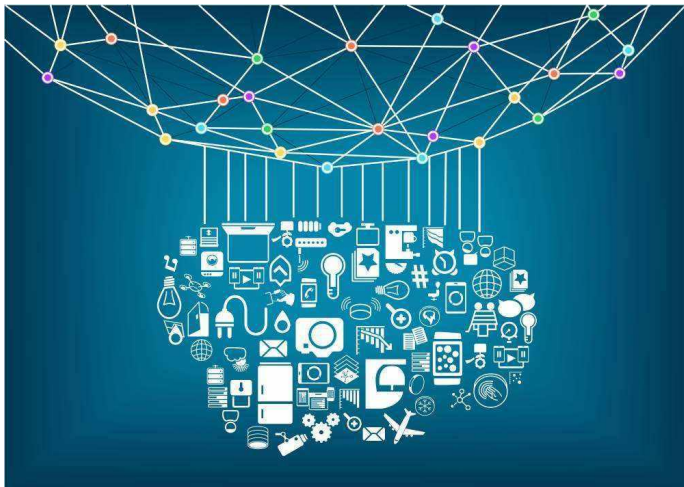
■ 自动机 形式逻辑 统计机器学习，汉语语言学 形式语法理论

➤ 经验主义（统计方法）：Statistics-Based Language Model

■ 线性：N-Gram; Bayes; Maximum Entropy; Hidden Markov Model
; Conditional Radom Fields; Support Vector Machine;

➤ 脑认知：神经网络，深度神经网络; CNN; RNN; LSTM;





大数据智能应用



NLPIR大数据搜索与挖掘实验室

定位：多语种·多模态·语义智能一流研究中心



20年NLP技术积累

- ✓ 中文信息处理最高奖**钱伟长一等奖**
(国内唯一分词方向)
- ✓ 算法：20+自主可控全链路NLP模块
20M 缓存边缘计算
数据：10GB语料库，20亿语义知识库
知识：10+行业先验知识积累

人、事、物、时空增强分析

- ✓ 多模态融合：NLP+OCR+语音+图文比对的**语义增强分析平台**
- ✓ 启动：<100份小样本冷启动
分析：**KGB知识图谱**关联分析
生成：报告智能生成

巡场订阅模式

- ✓ 国内：400+本地部署标杆客户
军工、中央网信办、公安部、国研中心、人行、建行、中电科、航天科工、国家电网、华为等
- ✓ 全球：40万用户验证，新闻集团、韩国RSN、意大利大使馆、新加坡南洋理工、日立等



张华平 博士

- ICTCLAS汉语分词创立者
- 创建并运营NLPIR大数据语义增强分析平台
- 北京理工大学副教授，大数据搜索与挖掘实验室主任
- 中国人工智能学会多语种智能信息处理专委会秘书长
- ✓ 中文信息处理领域最高奖：钱伟长中文信息处理一等奖（全国唯一“分词”方向）
- ✓ 新疆自治区科技进步
- ✓ 第一届ACL-SIGHAN国际汉语分词大赛
- ✓ 国家973汉语评测
- ✓ 中央网信办、中宣部、公安部等部委特聘技术顾问
- ✓ 国办电子政务总体组评审专家



钱伟长一等奖证书



新疆科技进步二等奖



张华平教授受CCTV采访解读苹果FBI揭秘大战。



研究方向布局

- 基础算法：深度学习、小样本学习、终身学习、提示学习；
- 语义增强分析
 - 自然语言处理：中英文、小语种；
 - 多模态处理：OCR、语音识别
 - 自然语言生成：自动摘要、特定领域GIS推理与生成
- 知识图谱
 - 知识图谱构建与推理应用；汽车营销
- 社交网络分析
 - 特定群体跟踪；
- 应用布局
 - 智能情报：舆情、军事情报、市长热线、汽车营销
 - 智能文书：档案处理、文档核查、辅助生成



承担的科研课题

- 面向互联网信息的司法舆情监测与分级预警技术研究及系统研发 2018YFC0832304 “公共安全风险防控与应急技术装备”重点专项（司法专题任务）
- 国家自然科学基金面上项目2项，语义主题与社交关系融合的特定群体发现关键技术研究 61772075
- 国家242信息安全计划十余项；BB文本校对
- ZF文本生成某课题
- GF创新特区：终身自主学习、配置安全、内容安全
- GF战略：人工智能、未来网络、未来电子、生物交叉



基于NLPIR引擎的语义增强分析平台



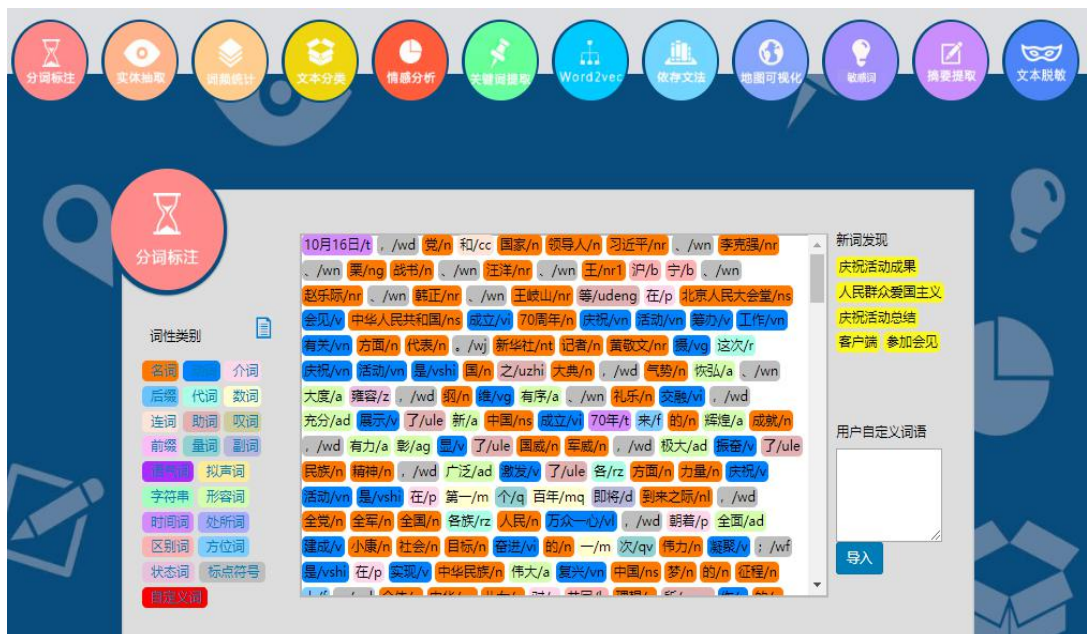
语义增强分析平台

- ✓ 具备70-80%应用功能
- ✓ 一周实现目标测试
- ✓ <100份样本冷启动
- ✓ 降低20-30%开发成本
- ✓ 提升系统智能化水平

核心技术：NLPIR语义分析引擎



“20年技术迭代
帮助全球40万用户实时理解中文”

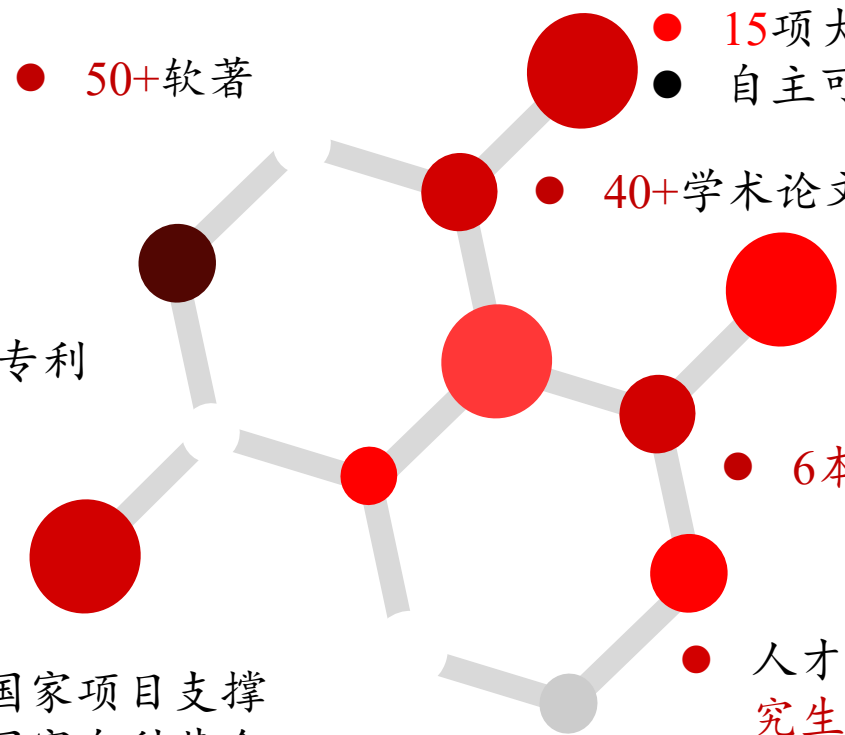


- 多语言智能分析：分词准确率接近99%，切分粒度可调整，融合20余部行业专有词典；新词发现、文本关键词提取：可处理海量网络文本数据（平均每小时处理至少50万篇文档）、关键词按照影响权重排序；一带一路小语种分析：维吾尔语、藏语、粤语、Hindi、缅甸、印尼
- 文本生成：自动摘要、报告生成
- 语义边缘计算：20M缓存即可部署
- Demo地址：<http://ictclas.nlpir.org/nlpir/>



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

主要成果

- 
- 50+软著
 - 3项专利
 - 10项国家项目支撑
 - 国家自然科学基金
 - 国家重点研发
 - 国家242
 - ZF预研/GF创新特区
 - 1个NLPIR大数据语义智能分析平台
 - 15项大数据语义分析算法组件
 - 自主可控，全面支持国产CPU/操作系统
 - 40+学术论文，SCI/EI 30+
 - 6本专著
 - 人才培养：70+博士硕士研究生，网信十佳讲师1人
 - 1000+论文引用
 - 服务在线用户11.9亿人
 - 近3年42万+技术开发用户
 - 直接经济效益12.4亿元



- **国家层面：**NLPIR已经应用于政府机构等国家部门，极大程度提高了其文本处理效率，提升了国家智能化水平。
- **社会层面：**NLPIR也出口到美国NCR、新闻集团、意大利ExpertSystem公司、韩国RSN公司、日本NEC与日立、新加坡南洋理工大学等国际知名机构。与华为、人民网、长安汽车也建立了合作；
- **科研层面：**NLPIR**无偿**提供在线分析平台和免费授权供广大自然语义处理的研究人员使用，就百度学术和知网对NLPIR/ICTCLAS的引用量已经均超过1000，Github的start量和Fork量也均超过1000。从2018起开发用户数超过**42万**。依据人民网2020年年报，人民网 PC 端、手机网、客户端、新媒体账号、代运营平台、聚合分发平台在全球覆盖直接用户数超过 **11.9 亿**。

奖励

科技



钱伟长中文信息技术处理科学技术一等奖



新疆维吾尔自治区科学技术进步奖二等奖、河北技术成果资质



新闻出版广电总局知识服务重点实验室-核心技术全国第二名



首席数据官联盟-中国大数据自然语言处理方向全国第一名



网络舆情挖掘系统、中文分词标注等20+项软件著作权



完美双数组管理和检索方法、互联网信息整合与更新方法等多项专利



学术认证



出版的专著

Huaping Zhang et, Big Data Intelligence (Book),
Tsinghua Univ. Press, 2019

Huaping Zhang et, Big Data Big Talk (Book),
Publishing House of Electronics Industry, 2019

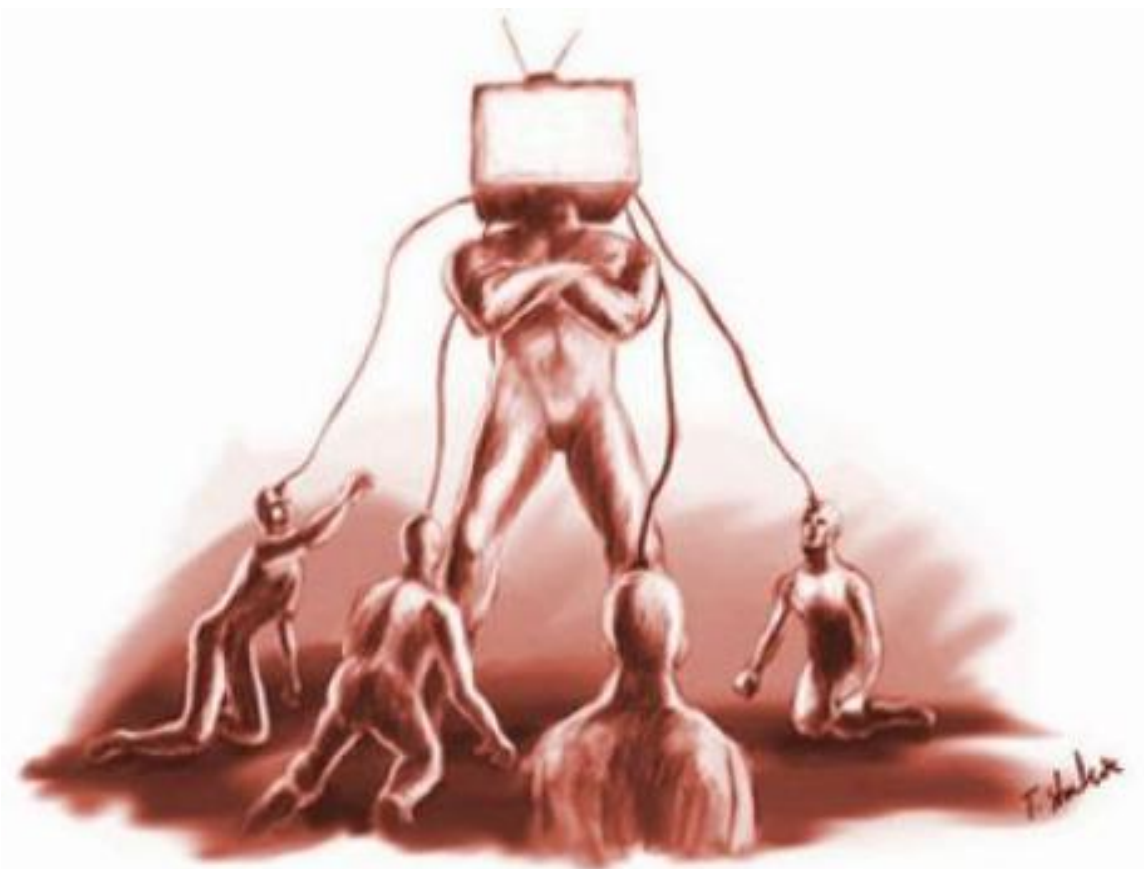
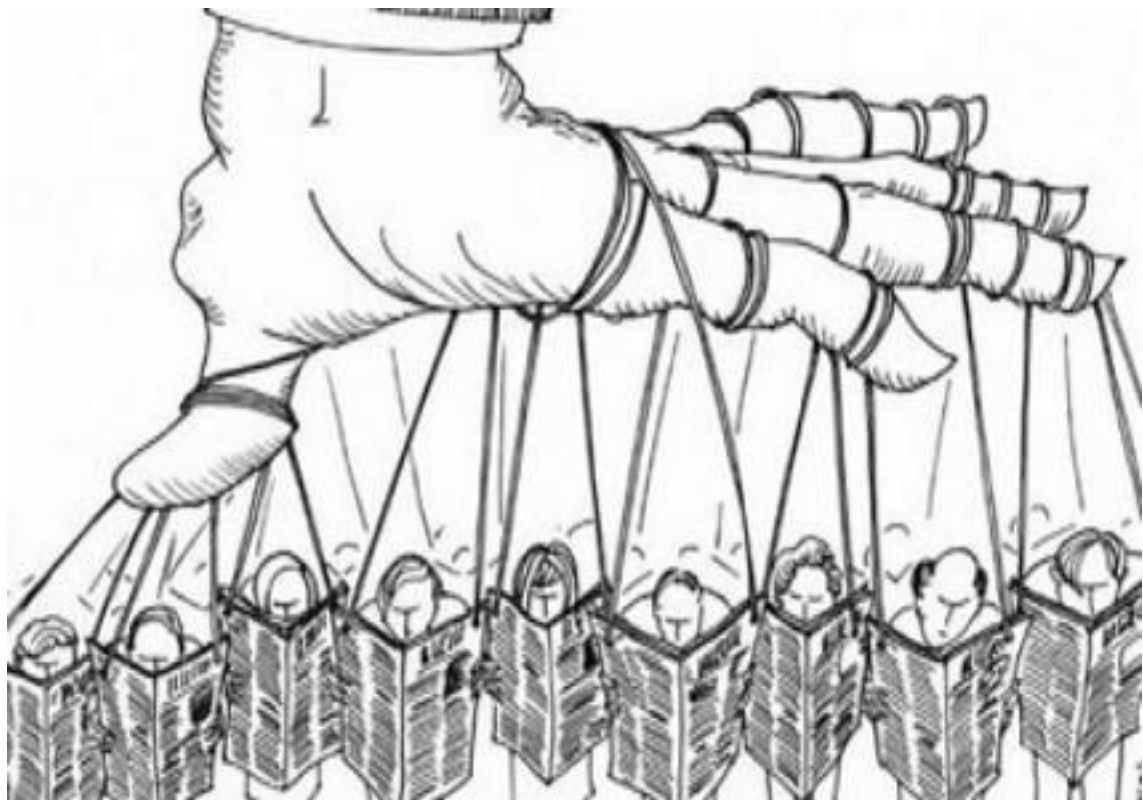
ISBN:978-7-121-30181-0

Huaping Zhang et, Big Data Search and Mining (Book),
Science Press, 2014.5 (ISBN:978-7-03-040318-6)





最恐怖黑客：黑掉你的脑子



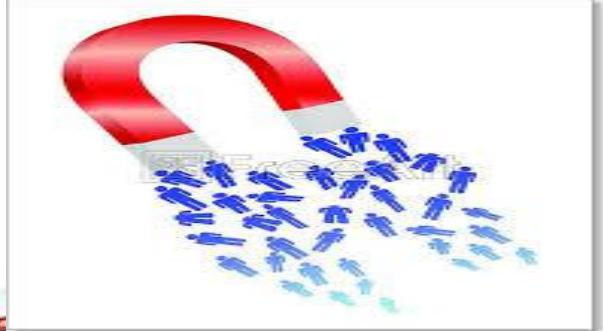
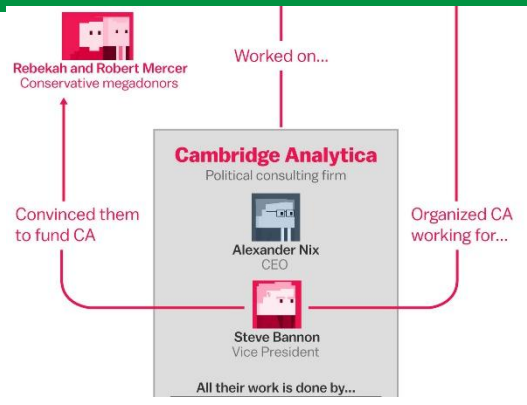


社交舆情操控

➔ 剑桥分析控制62国选举

➔ 美国2016大选社交网络干预

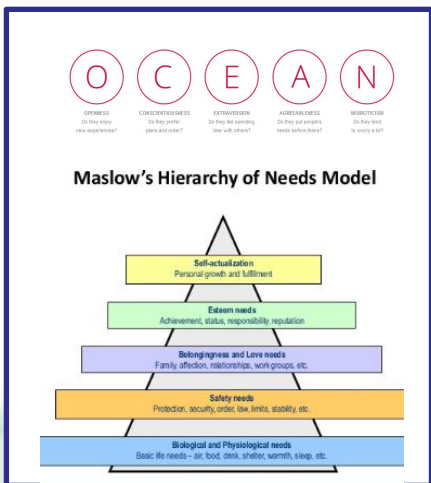
➔ ISIS全球吸引恐怖分子





舆论操控者

心理学模型



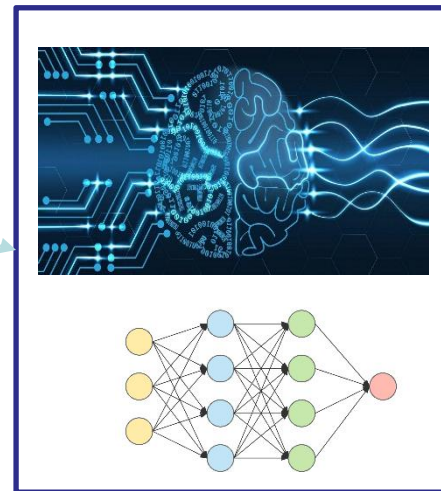
目标群体

社交舆情操控

社交网络



计算模型



舆论操控



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY



剑桥分析:方法论

5,000 data points per person

We collect up to 5,000 data points on over 220 million Americans, and use more than 100 data variables to model target audience groups and predict the behavior of like-minded people.

Constant testing & improving

Our data scientists and psychologists are constantly testing new modeling and research techniques to ensure all our data sets and audience segments are the most advanced in the market.

OCEAN and the Big Five

We use the established scientific OCEAN scale of personality traits to understand what people care about, why they behave the way they do, and what really drives their decision making.



Spot the differences

Understanding the complex web of OCEAN personality traits behind behavior lets us see why people who look similar on the surface often want and respond to completely different things.



北京理工大学
INSTITUTE OF TECHNOLOGY



社交輿情操控

Conversations that move people

When you go beneath the surface and learn what people really care about you can create fully integrated engagement strategies that connect with every person at the individual level.

Same demographics, different personalities



Female
25-35 Years old
AMEX User



People with high openness and extraversion love new experiences they can share with lots of people.



Female
25-35 Years old
AMEX User



People with low openness and extraversion really value down time spent with their closest friends.

We Call This Behavioral Microtargeting

Discover. Understand. Engage. Repeat.

Combine our full suite of data-driven audience insight and engagement techniques with our unique and powerful Behavioral Microtargeting service that constantly learns, improves and delivers.

With Behavioral Microtargeting you'll be able to anticipate the needs of your customers and predict how their behavior will change over time, so you can build services, products and campaigns they really love.



Geographic View



Demographic View



Psychographic View



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY



ISIS: 网络极端势力如何吸引恐怖分子

Table 1 Offender behaviour characteristics

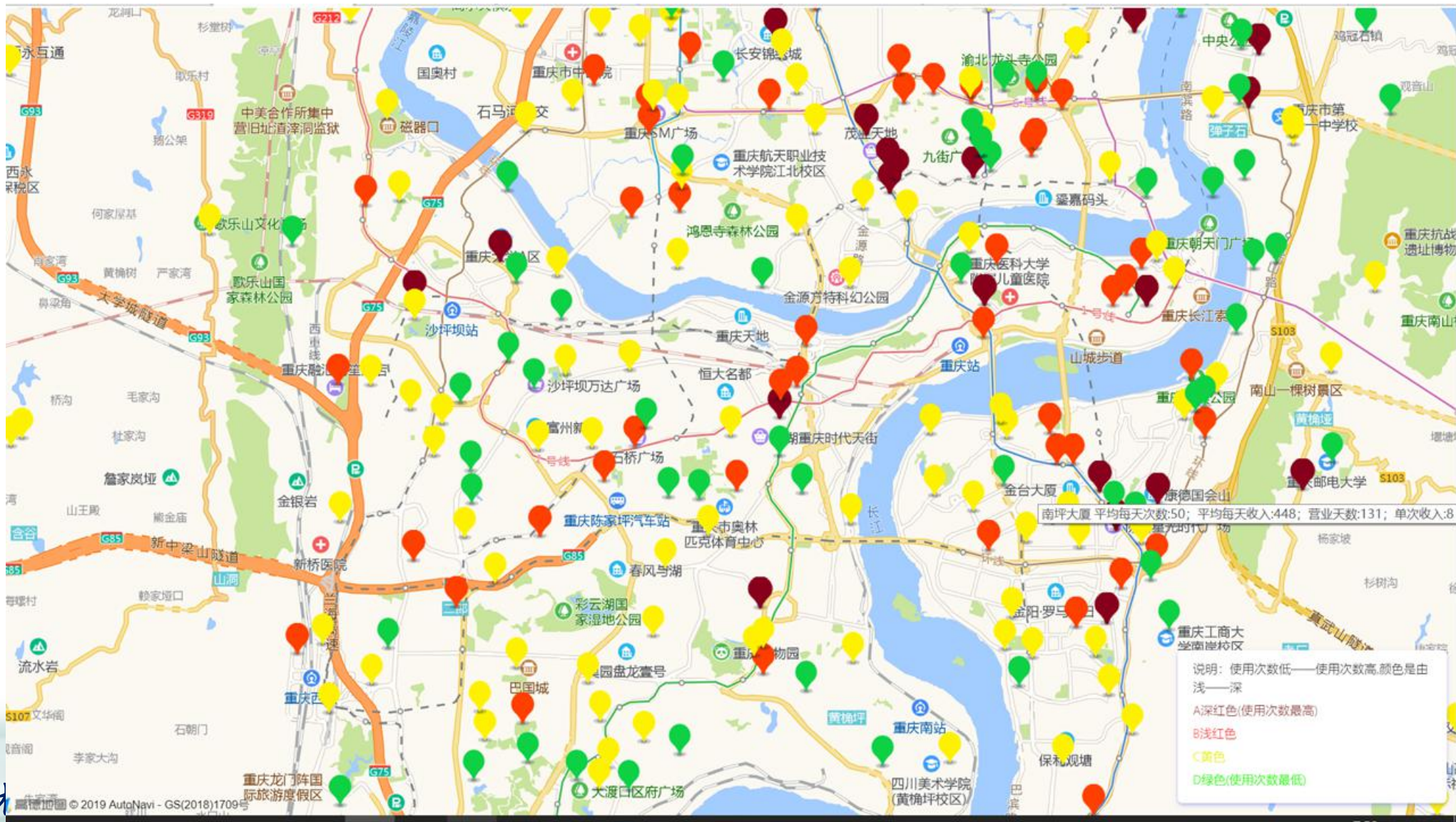
Type	Characteristics	Cases on Twitter	Cases on Facebook	Total No of Cases
Cyber Mobs	Using social media platforms to create a mob mentality and urging others to fight for the Isis goal. This is done through group posts, videos and comments of hate directing groups of Muslim's to fight. Often personified through retweets, likes and views of specific Isis propaganda materials.	78	55	133
Loners	Often done through individual posts and comments. This individual is someone who is attracted to the Isis campaign but clearly is exposed to individual grievances and has a lone mentality.	51	65	116
Fantasists	Someone using social media platforms to fantasise over the Isis movement. In particular, these individuals have blurred the lines between reality and fiction and are making direct plea's to fight for Isis.	45	94	139
Thrill Seekers	People who are promoting Isis propaganda through videos and posts and forums. Indeed, some of these individuals claim to be directly using the Internet for online extremist purposes. These individuals are describing the sense of adrenaline rush they are receiving by watching and partaking in fighting on the battlefield whether online or offline.	85	98	183
Moral Crusaders	These individuals are talking about the moral duty to fight. Many of these individuals are also constructing arguments based on ideology and theology as a means to promise people external rewards.	140	95	235
Narcissists	These people are using political, foreign policy and individual grievances as a means to whip up a climate of revenge seeking and wanting to fight for the Isis mission and goals.	166	104	270
Identity Seekers	Mostly this is users who appear to be seeking some form of identity. Primarily people searching for some form of masculinity and therefore the Isis recruitment drive appeals to them. This applies to males and females.	87	101	188



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

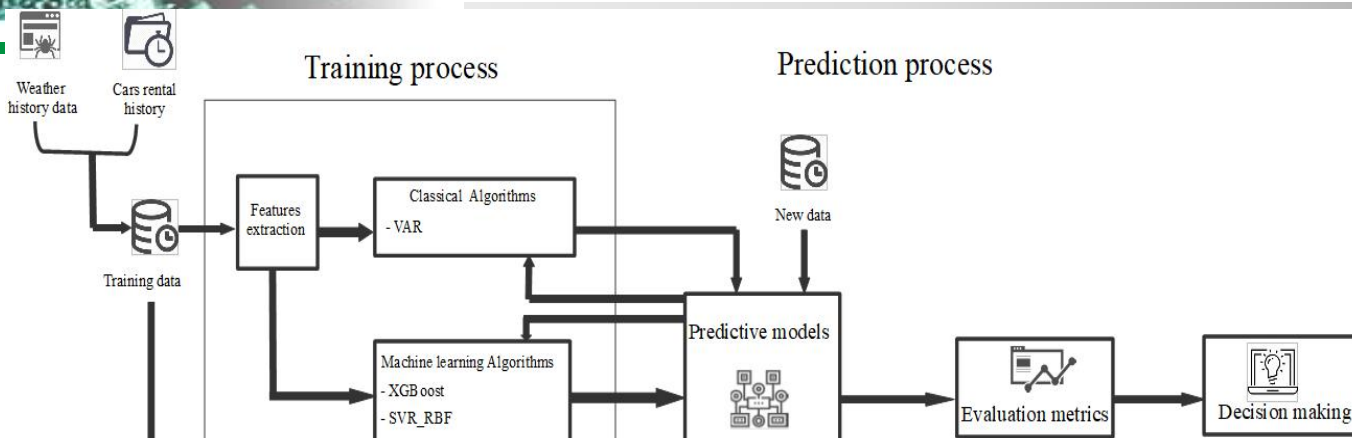


某汽车出行宏观画像





某汽车出行租车预测模型



	History + weather + weekend/weekday			History + weather			History + weekend/weekday			History		
	MAE	MSE	RMSE	MAE	MSE	RMSE	MAE	MSE	RMSE	MAE	MSE	RMSE
Class A (#16)	0.29858	0.35530	0.59607	0.40409	0.51835	0.71996	0.44704	0.61301	0.78295	0.68891	0.68891	0.83000
Class B (#104)	0.15430	0.17264	0.41550	0.21268	0.23436	0.48411	0.30067	0.38240	0.61839	0.59425	0.59425	0.77087
Class C (#6)	0.25730	0.31151	0.55813	0.32569	0.49750	0.70534	0.42327	0.72102	0.84913	0.78899	0.78899	0.88825
Class D (#28)	0.00083	0.00083	0.02888	0.00250	0.00250	0.05002	0.02335	0.02335	0.15282	0.06839	0.06839	0.26152
Class E (# 25)	0.01168	0.01251	0.11185	0.01418	0.01501	0.12253	0.01585	0.02419	0.15552	0.02419	0.02419	0.15552

Table 7: evaluation metrics for multivariate time series forecasting using LSTM

	History + weather + weekend/weekday			History + weather			History + weekend/weekday			History		
	MAE	MSE	RMSE	MAE	MSE	RMSE	MAE	MSE	RMSE	MAE	MSE	RMSE
Class A (# 16)	0.90898	2.02258	1.42217	0.90889	2.02560	1.42323	1.50311	3.15878	1.77729	1.50419	3.16149	1.77805
Class B (#104)	0.97605	2.11016	1.45264	0.97571	2.11519	1.45436	0.97568	2.11486	1.45436	0.97568	2.11487	1.45259
Class C (#6)	0.90898	2.02258	1.42217	0.90889	2.02560	1.42323	0.90901	1.42194	2.02192	0.90884	2.02483	1.42297
Class D (#28)	0.60189	0.66971	0.81836	0.60221	0.66972	0.81836	0.60192	0.81836	0.66971	0.60225	0.66972	0.81836
Class E (#25)	0.16379	0.10897	0.33011	0.16437	0.10903	0.33013	0.163867	0.10898	0.33013	0.16449	0.10905	0.33023

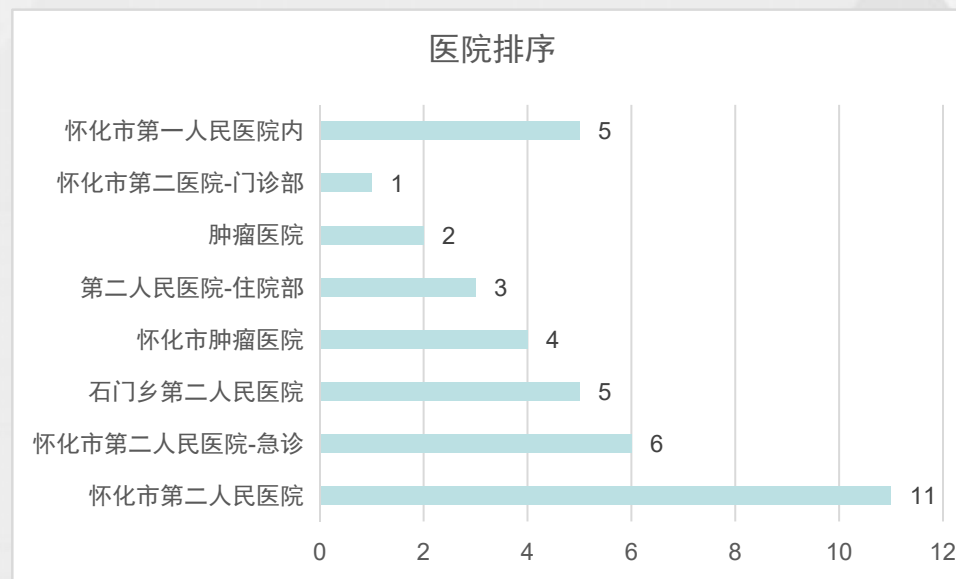
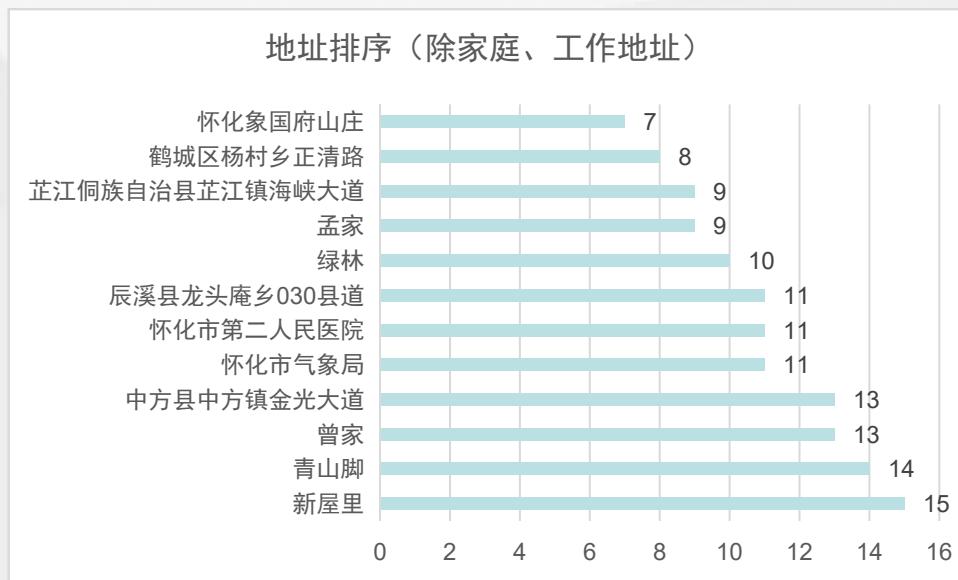




GPS人物画像：POI

➤ 除去可能是家庭、工作的地址，选其它高频地址分析4人除家和公司常去的地方。

湖南怀化气象工作人员



医疗：怀化市第二人民医院共出现32次，此人经常开车去医院，属于紧急情况用车较多的类型。**其他高频地址比如：**新屋里、青山脚、曾家、绿林、孟家都是村庄，距离家庭住址盛世嘉园车程约一小时左右，不属于休闲娱乐场所，可能是探亲访友用车。

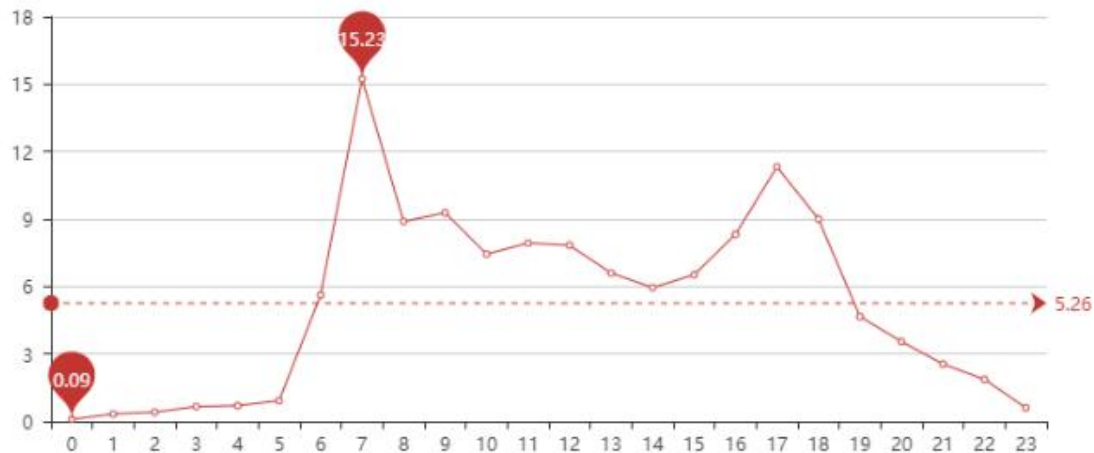


GPS人物画像：工作性质

工作日每天各时段开车时长=工作日总开车时长/工作日总天数

湖南怀化气象工作人员

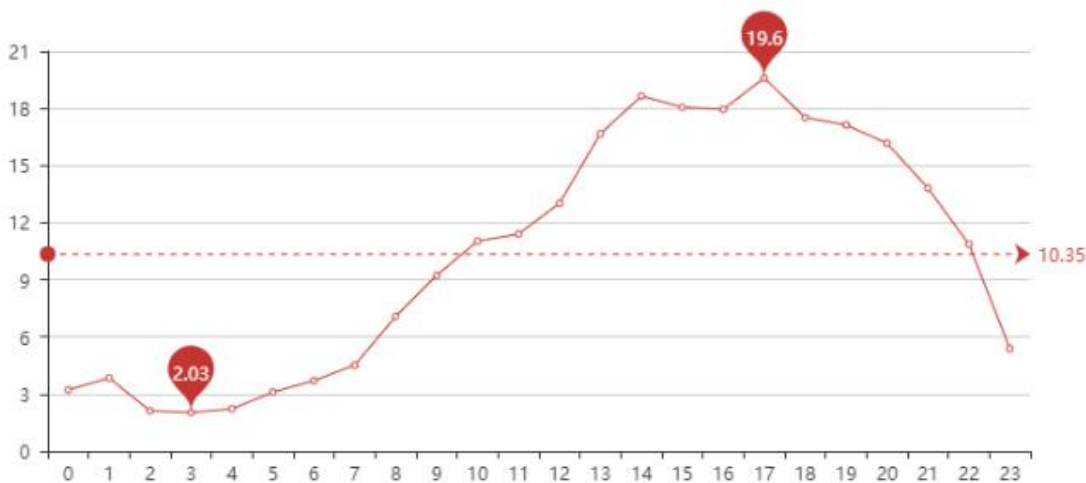
○ 工作日每天各时段开车时长



高峰时段：7点和17点，典型早8晚5型用车
工作性质：政府单位

广东云浮城镇小青年

○ 工作日每天各时段开车时长



午后用车型：用车高峰时段14-22点
工作性质：非正式单位，自由职业者

智能机器阅读



摘要提取



文档支持

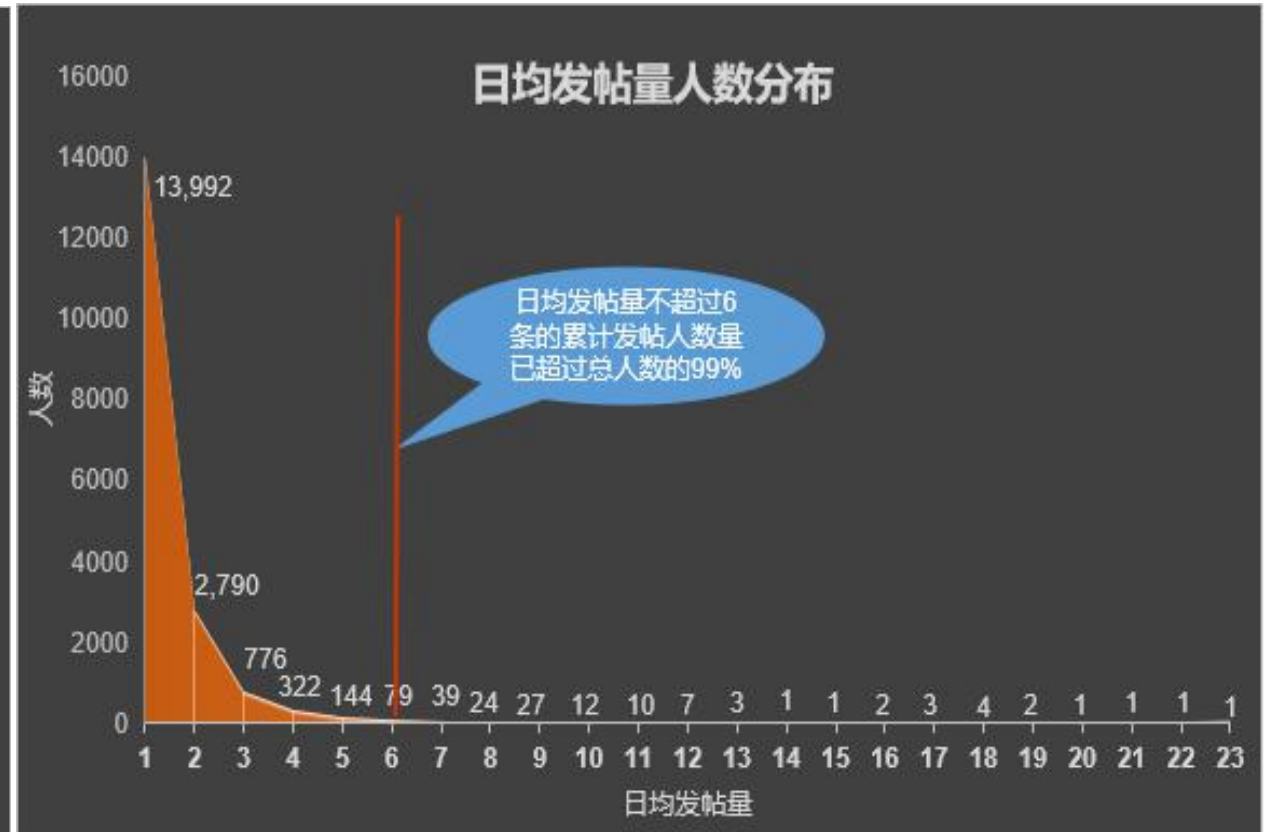
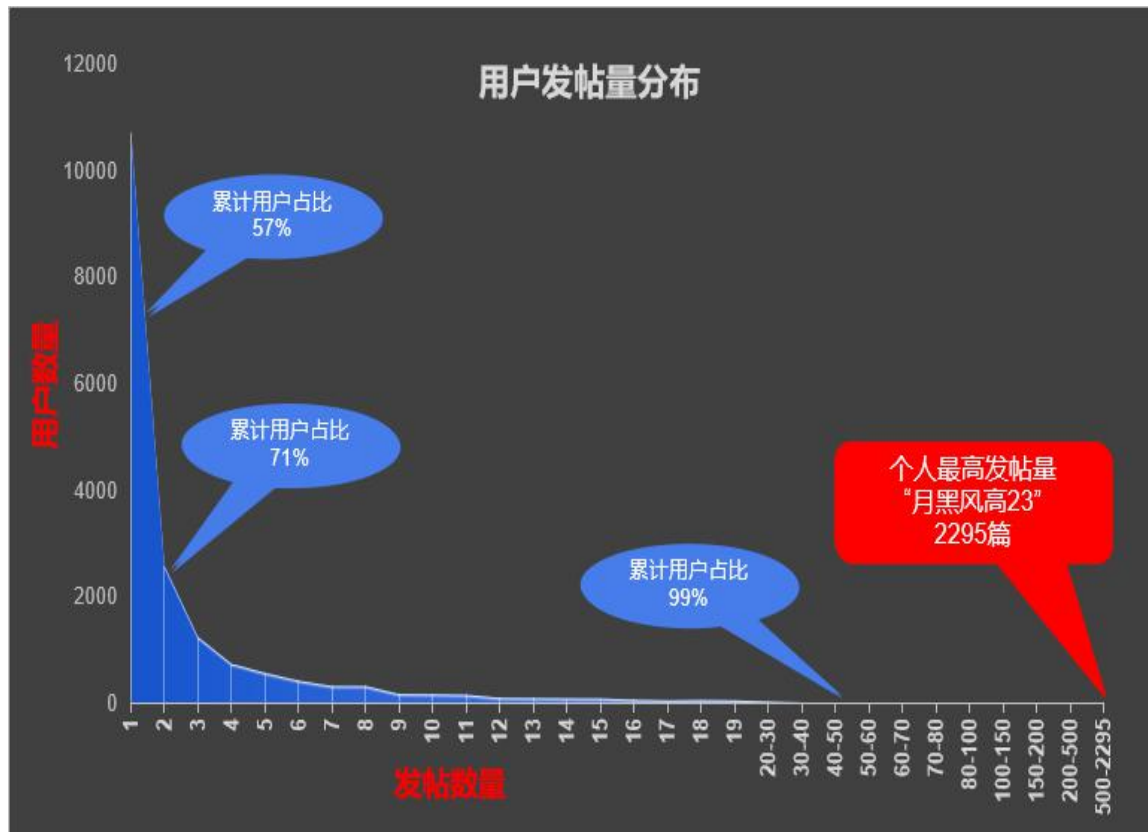
东盟加快发展数字经济（国际视点）--国际--人民网本报驻泰国记者刘慧2021年08月11日05:21来源：人民网 - 人民日报
盟地区多国近期持续出台数字化发展战略，促进数字经济加速发展。数字支付业务也在东盟国家呈上升态势。中国和东盟国家抓住数字化转型机遇，在数字基础设施、5G、人工智能、大数据等领域打造更多合作亮点。

繁
拼

东
本
20
东
中
东
西
数
在
票
这
疫
今
送
在
杯

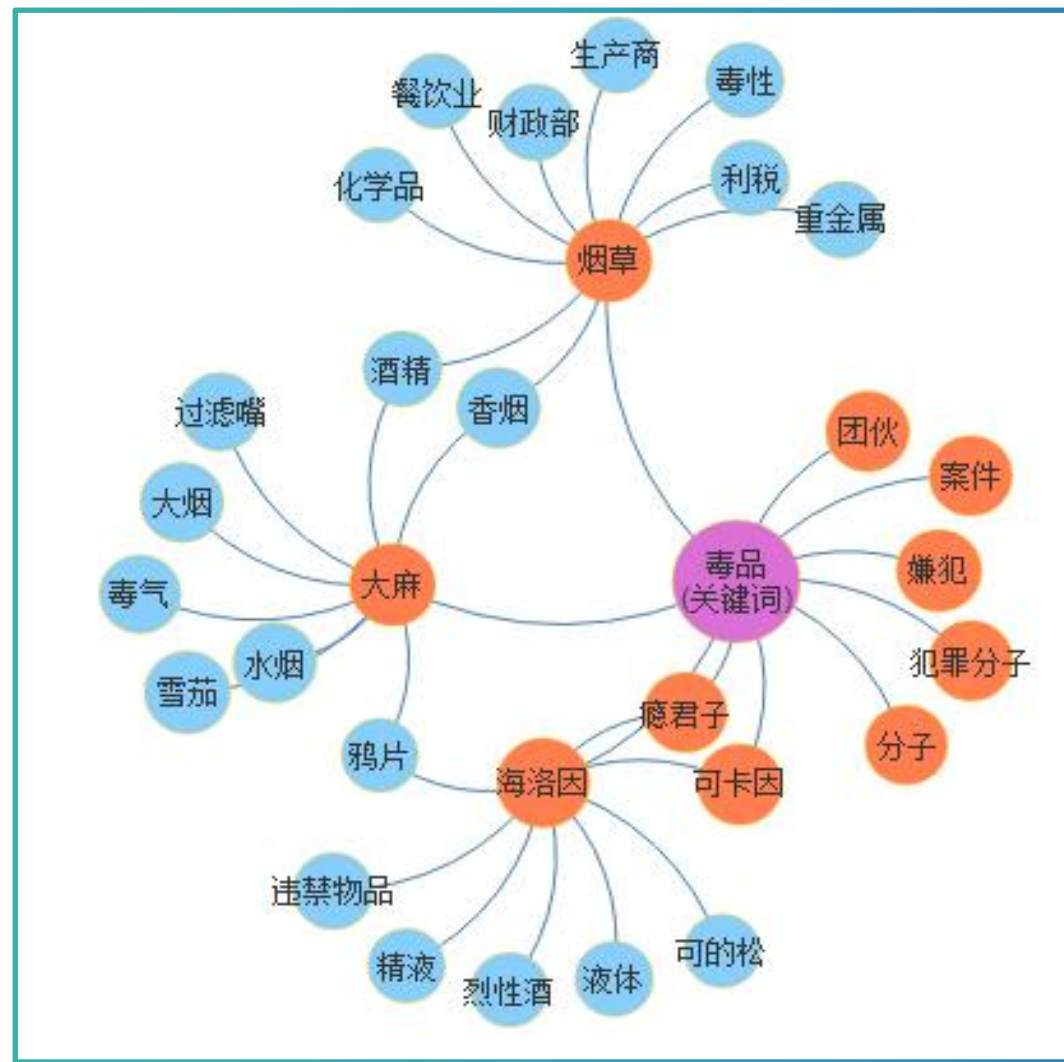
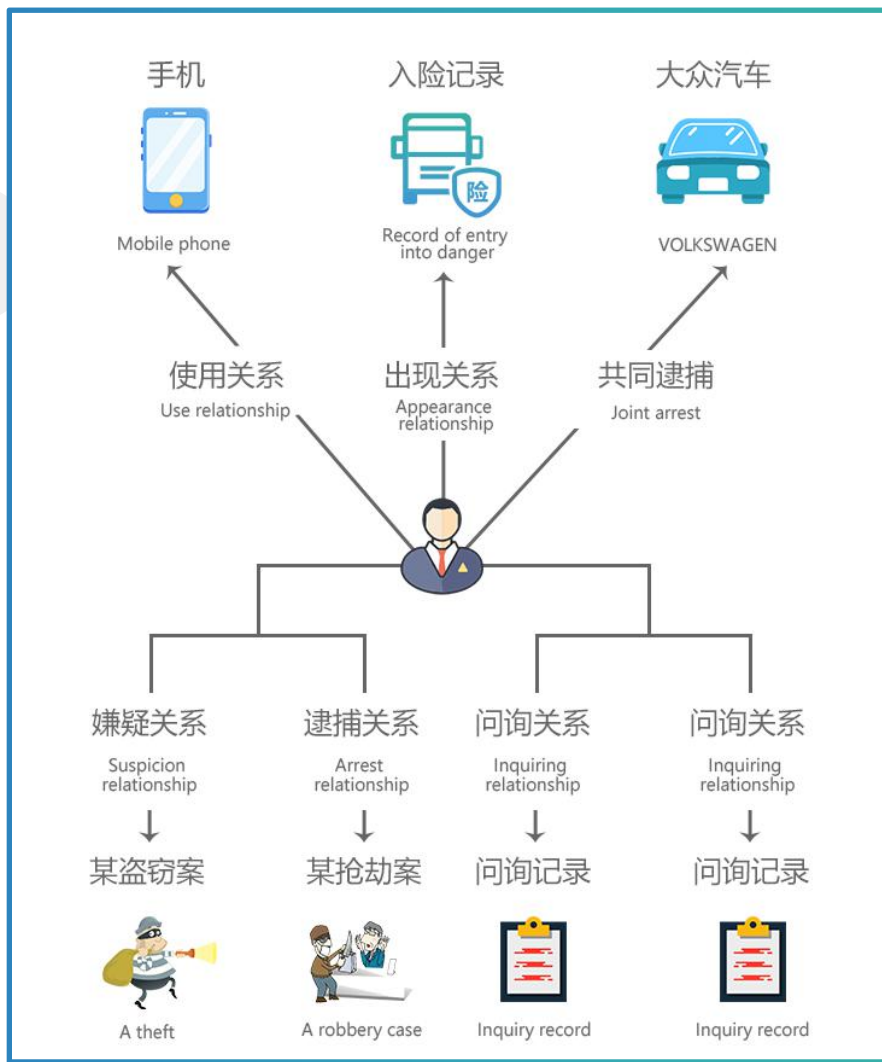
某品牌舆论场画像

发帖人总量：**18645个**； 发帖总量：**87484条**； 人均发帖量：**4.7条/人**



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

公安语义增强应用

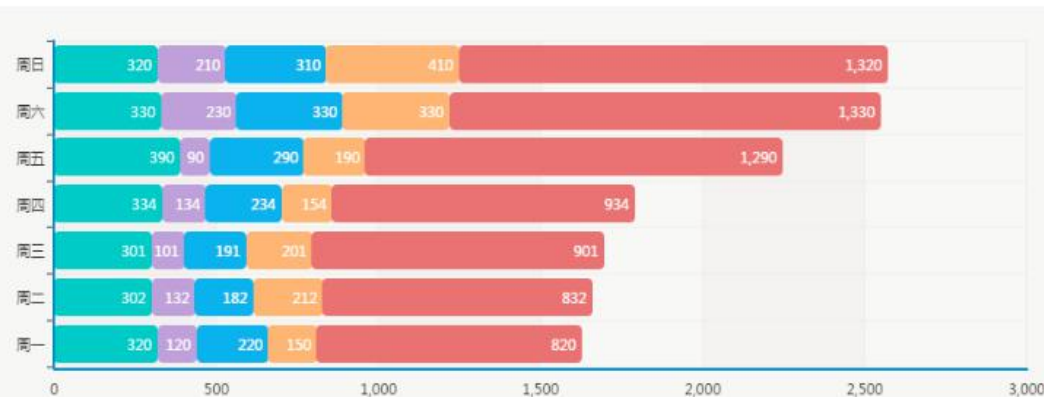


公安语义增强应用

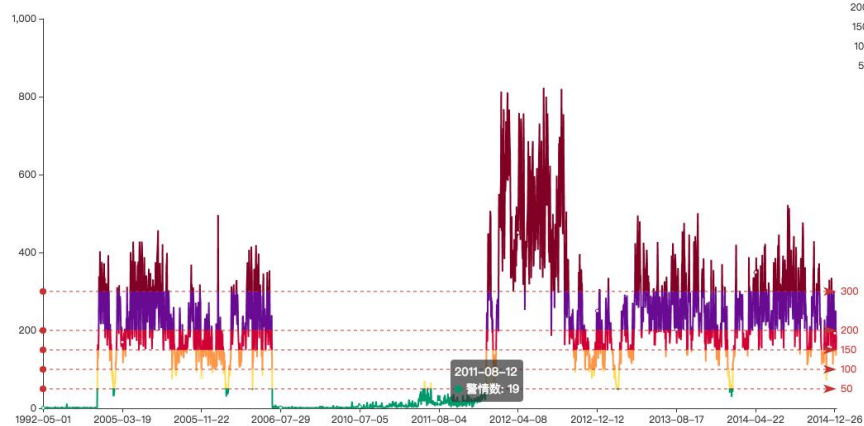
聚类高发

- Cluster 0: 天作 使用费 非法占有
- Cluster 1: 伪造 交通事故 赔偿金 保险公司
- Cluster 2: 博亿 房屋买卖 虚构
- Cluster 3: 鲁迅文学奖 作家协会 诗词
- Cluster 4: 骗走 人民币 平房 公墓 灵泉 塔位 灵塔 地葬位
- Cluster 5: 押金 天纺 柯东 冒用 丰体
- Cluster 6: 为名 培训 假借 团伙 考证 帮助

现实高发



西城区警情历史数据



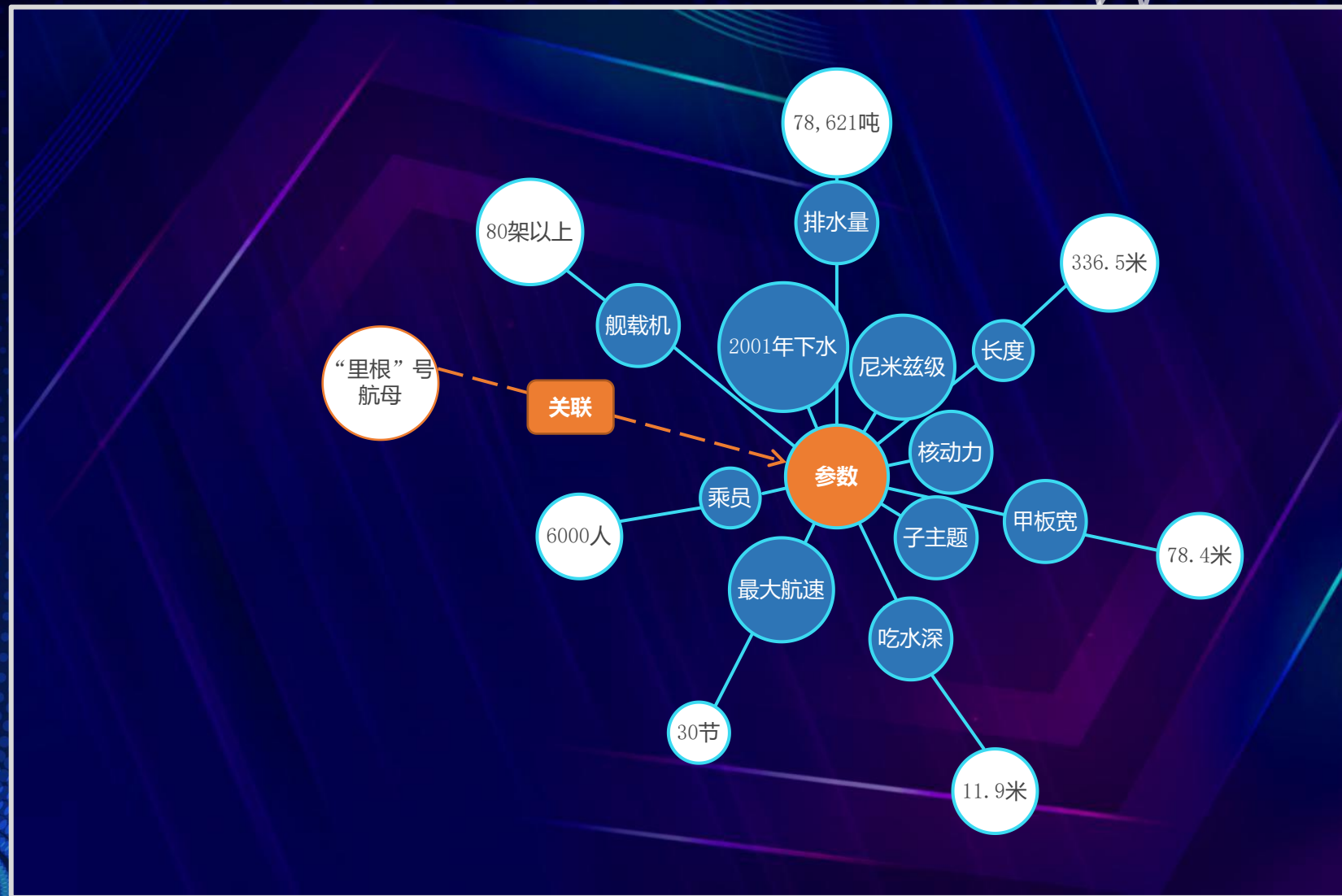
预警提示



军事知识图谱：从浅层信息挖掘深度情报



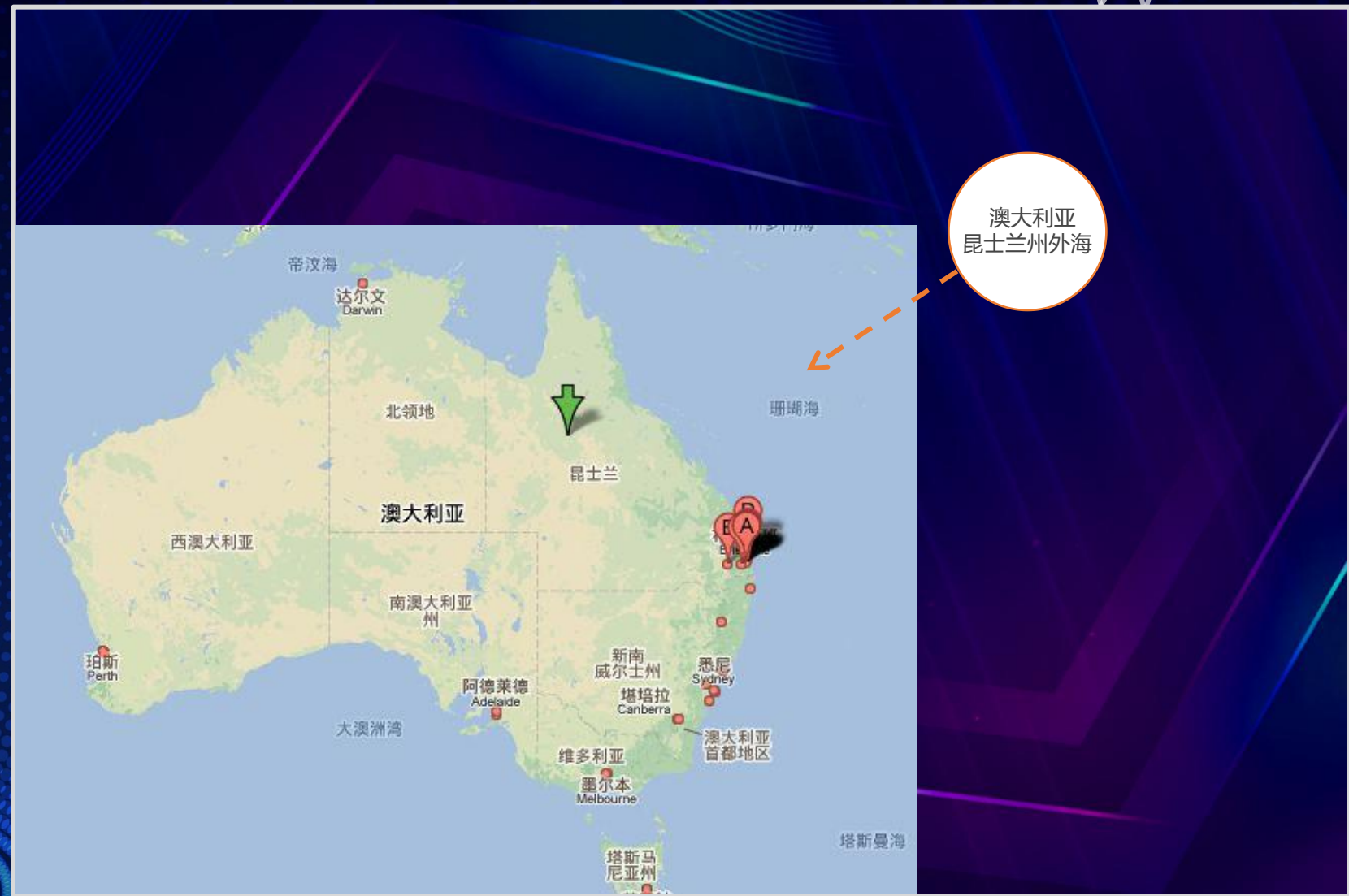
军事知识图谱：从浅层信息挖掘深度情报



军事知识图谱：从浅层信息挖掘深度情报



军事知识图谱：从浅层信息挖掘深度情报



智能机器阅读：篇章结构理解

2006年以来，樟村坪镇委、镇政府在市委、市政府，区委、区政府的正确领导下，坚持以邓小平理论、“三个代表”重要思想和党的十六大精神为指导，全面落实科学发展观，转变发展观念、创新发展思路，坚持走矿山立镇、项目强镇、种养富镇、和谐兴镇之路，紧紧围绕发展第一要务，突出资源保护和综合开发利用，加强企业管理和项目建设，加大基础设施建设和农业产业结构调整，扎实推进全镇各项工作，全力打造湖北磷化第一镇。

2006年1月至5月，全镇完成企业总产值6.56亿元，比去年同期增长44.08%；完成工业总产值3.89亿元，同比增长31.31%；实现财政收入5441万元，同比增长121%，全年财政收入有望突破亿元大关。

一、以矿山管理为重点，培育经济核心竞争力

2006年，全镇矿业经济立足市场促规范，提高效益保安全，稳步推进工业化，在规范、安全、效益、竞争力、保值增值方面有所突破。工业产值增长15%以上，规模工业产值增长20%以上，安全生产力争实现零事故目标。

1. 主动参与，积极配合，进一步健全矿山管理的长效机制。

"docName": "情况汇报材料.docx",

"childs": [

 {"Text": "中共樟村坪镇委樟村坪镇人民政府工作情况汇报材料\r",

 "Style": @Object{...},

 "level": 0,

 "childs": [

 @Object{...},

 @Object{...},

 {

 "Text": "以矿山管理为重点，培育经济核心竞争力 \r",

 "Style": @Object{...},

 "level": 1,

 "childs": [

 @Object{...},

 {

 "Text": "主动参与，积极配合，进一步健全矿山管理的长效机制。 \r",

 "Style": @Object{...},

 "level": 2,

 "childs": [

 {

结构抽取

可点击key和value值进行编



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

智能辅助写作：文书自动核查

工程名称：小学项目总承包工程

采购供应 一、材料名称、规格、数量及价格（见后附表）

序号	材料名称	规格型号	单位	数量(暂估)	含税综合单价(元/个)	含税总价			
1	潜水泵	D	序号	材料名称	规格型号	单位	数量(暂估)	含税综合单价(元/个)	含税总价
2	污水泵	2.	1	潜水泵	DN40	台	1	2900	2900
3	污水泵	3.	2	污水泵	2.2KW	台	2	850	1800 1700.00 [ErrorMsg:1000:数量价格金额核算]
4	塑料管	D	3	污水泵	3.0KW	台	2	950	1900
5	排水管	D	4	塑料管	DN40	米	100	8	800
合计			5	排水管	DN65	米	60	7.5	450
合计						7750			
						发票税率：3%			

核数总

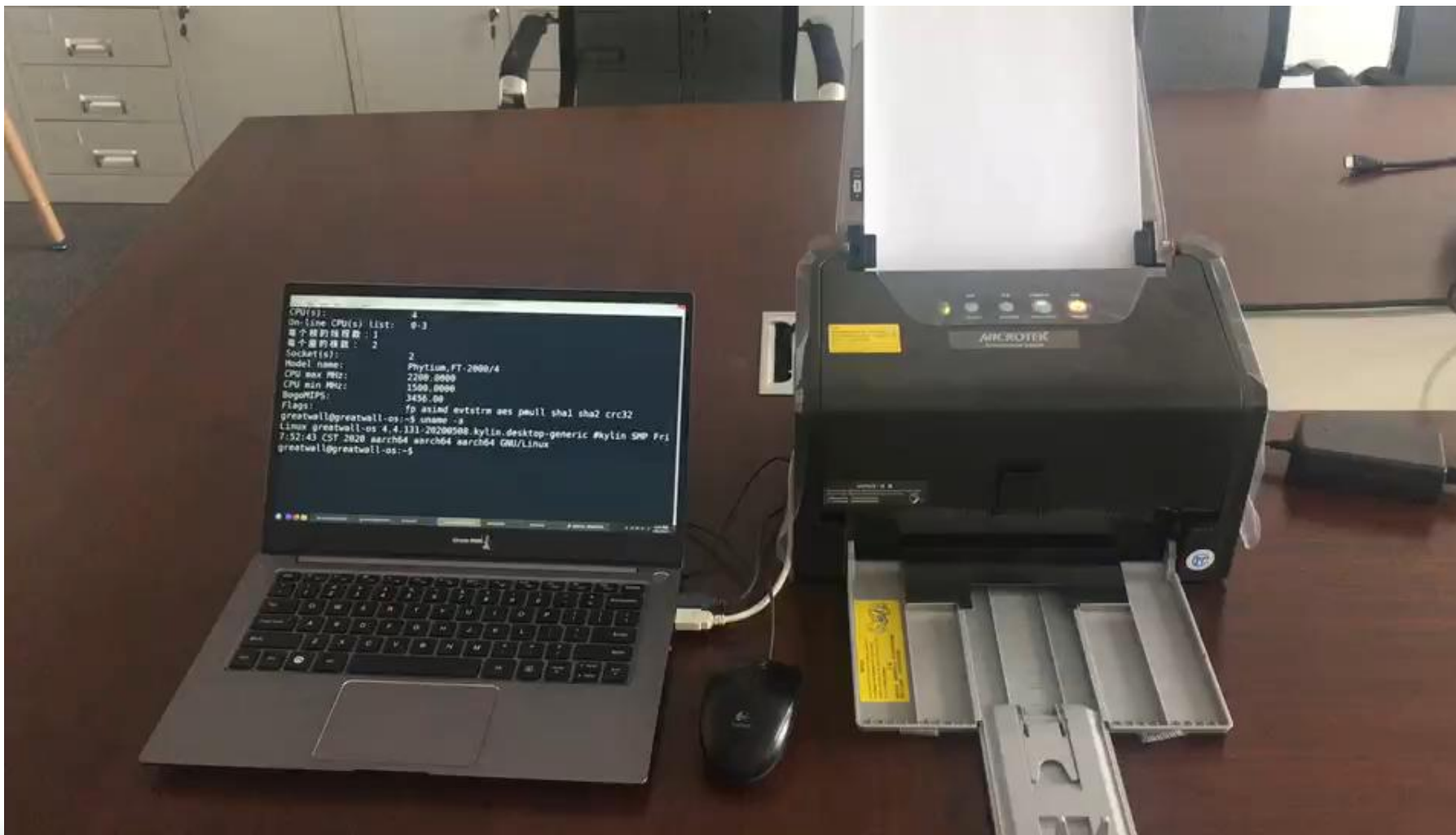
金额核算核查：
金额大小写不一致；
分项总价之和不等于
合计总额。

1、本合同价款为：~~7751~~7750
[ErrorMsg:金额大小写不一致]7850.00
[ErrorMsg:1000:分项总和必须等于总额]元，
(大写)：柒仟柒佰伍拾元整；

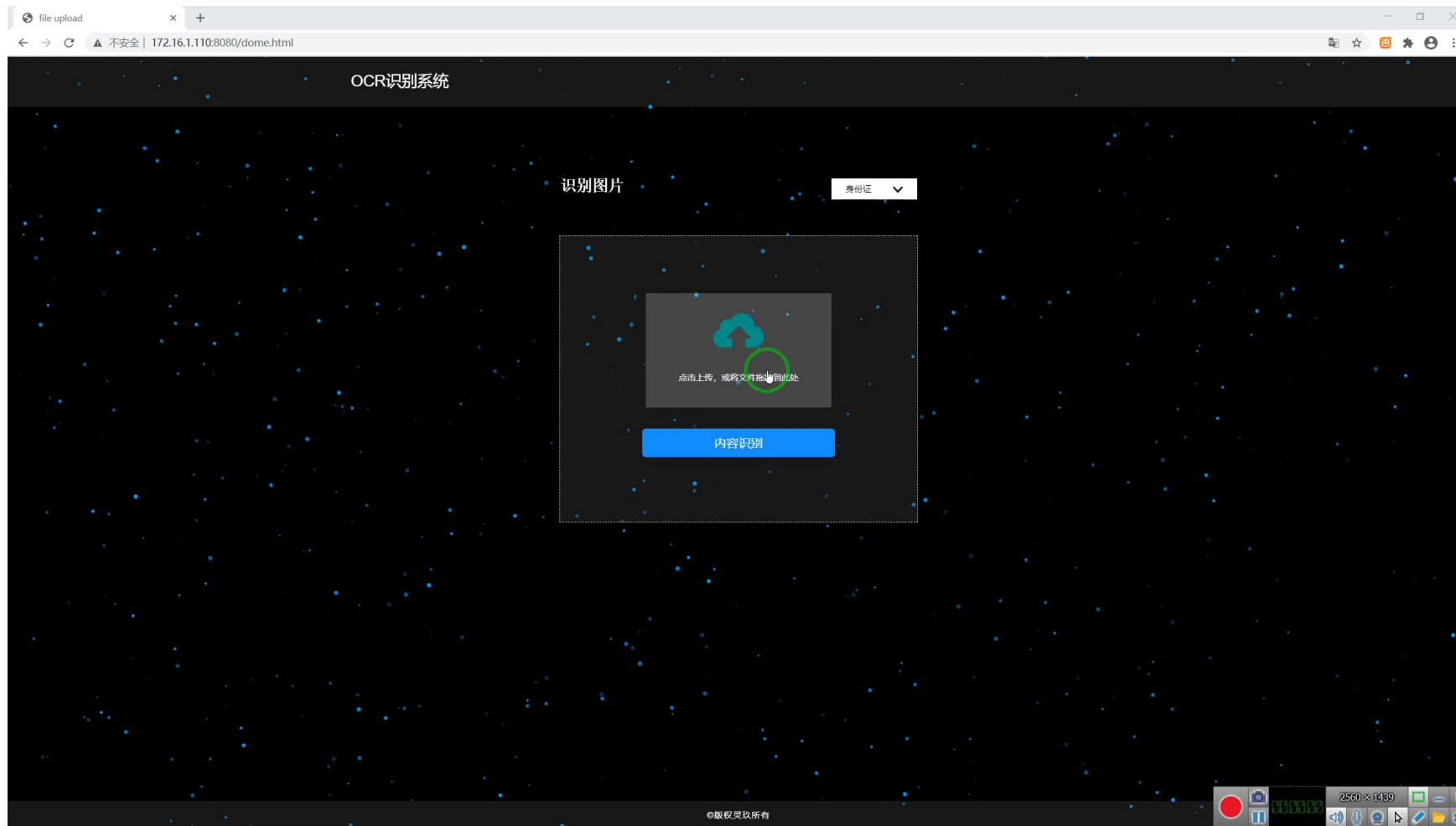
2、合同价款采用固定单价方式确定，合同单价包括材料费、包装费、运输费、仓储费、利润、税金等材料出厂运至施工现场指定位置的所有费用，结算时不再调整合同单价；

采购方：（盖章）
[ErrorMsg:实体前屏
datetime="2018-09-19T1
中黑名单]（盖章）
采购方：中航天建设
供应方：固安轴承
工程名称：小学项目总
施工地点：固安县

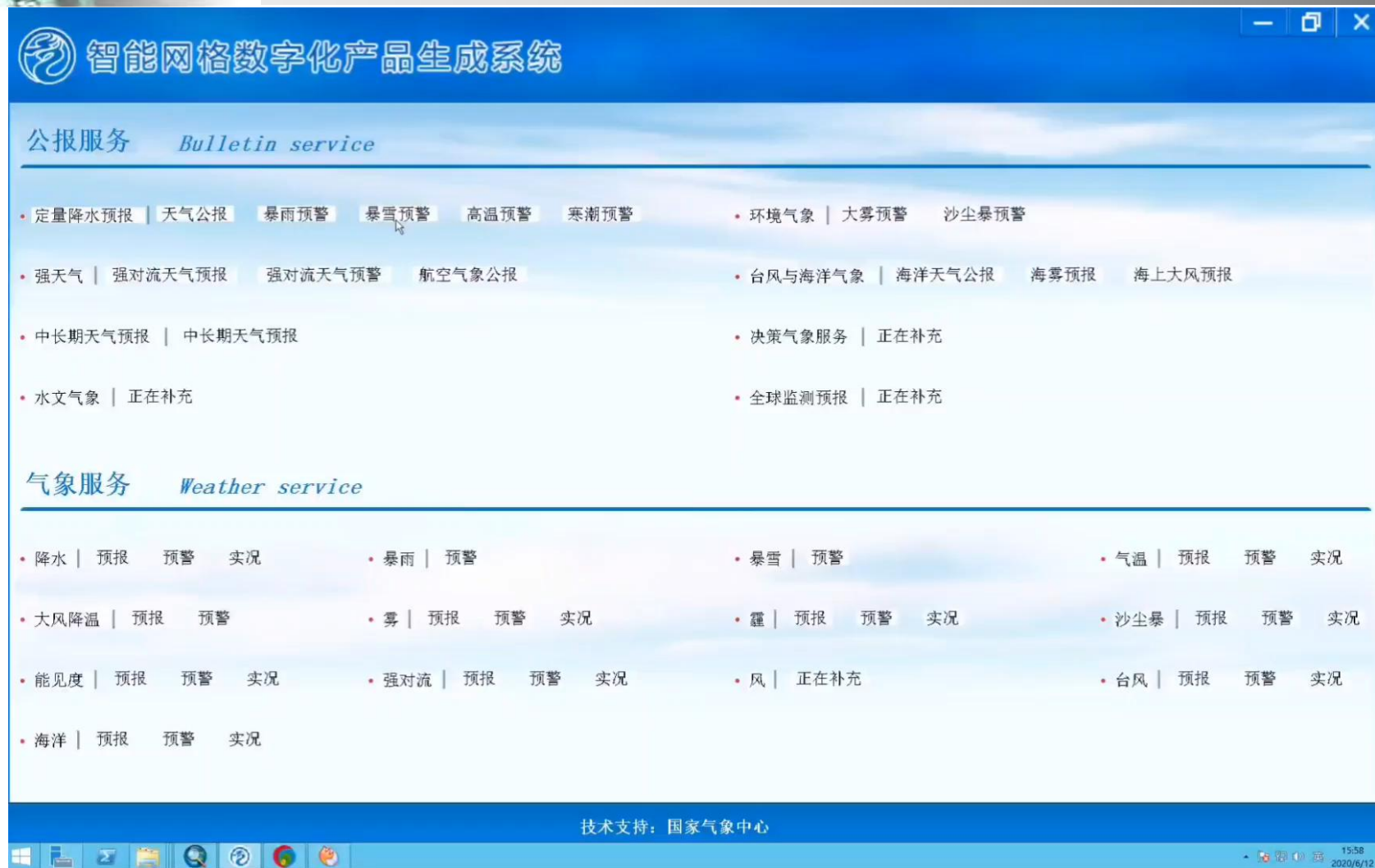
智能辅助写作：批量国产化扫描识别



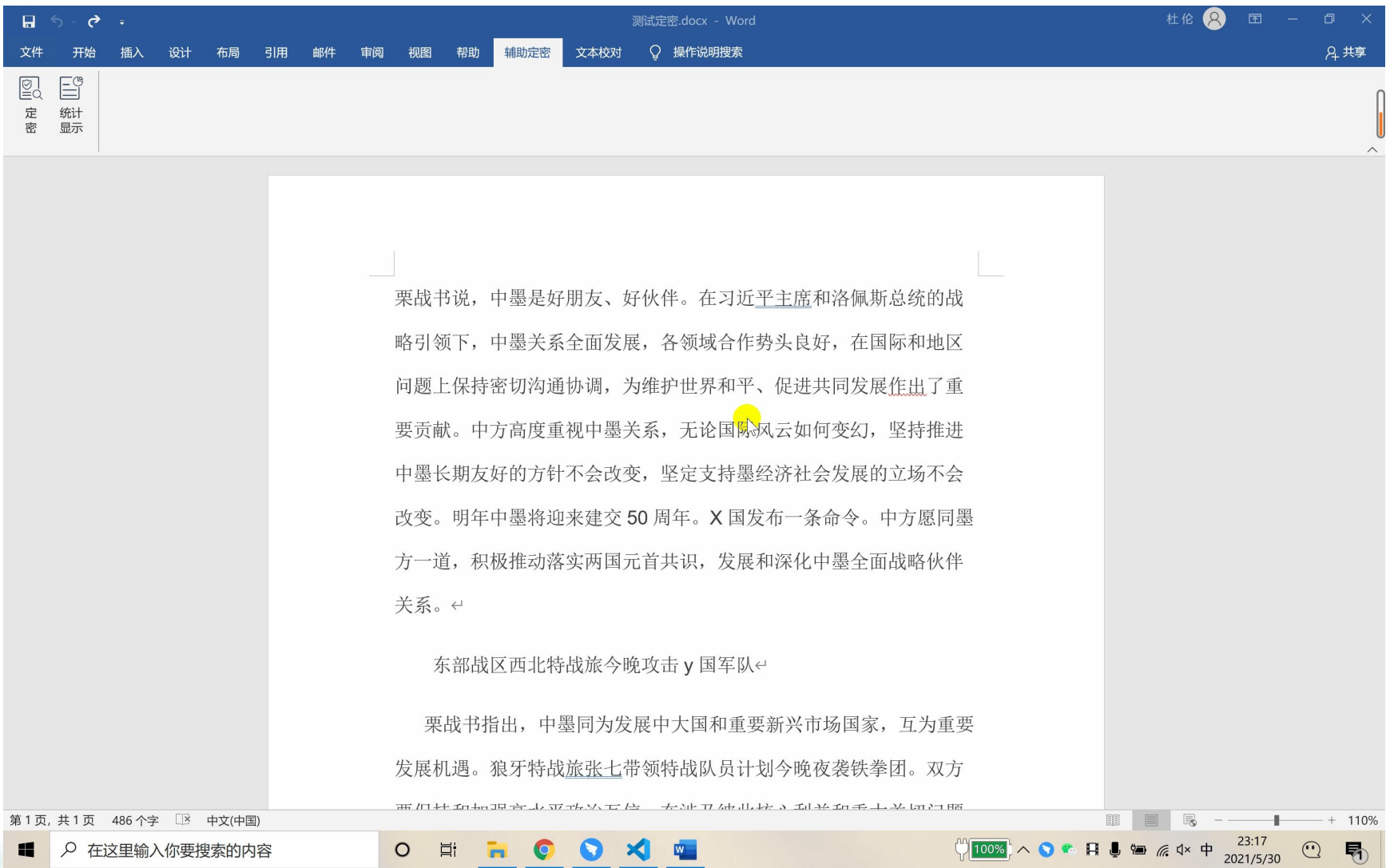
智能辅助写作：图文表识别



智能文书写作：气象文本生成



智能增值服务：辅助定密



多模态分析



- **人工智能**：现代科学皇冠上的明珠；
- **自然语言处理**：**人工智能**皇冠上的明珠
- 机器一思考，人类就发笑；人类一思考，上帝就发笑。





感谢关注聆听！



张华平

Email: kevinzhang@bit.edu.cn

微博: @ICTCLAS张华平博士

实验室官网:

<http://www.nlpir.org>



大数据千人会

