

大数据分析实践

Introduction

Qiong Zeng (曾琼)

qiong.zn@sdu.edu.cn



演练画面



勝利日
1945-2025 胜利日将举行盛大阅兵
天安门阅兵是纪念活动重要组成部分

1945-2025

1945-2025

阅兵式里的大数据分析系统

- **数字力量**：基于大数据分析的毫米级精度方阵（实现多单位协同指挥与精准行动）、北斗定位系统（提供实时精准定位与数据支撑）、量子通信车（保障数据传输效率与安全）、东风-51导弹（依托数据链实现智能目标锁定与打击）、智能防空体系（通过大数据实时处理威胁信息）、无人机集群控制系统（基于算法协同执行复杂任务）。
- **装备安全保障**：高空高速无人机可依托数据分析精准识别与压制航母战斗群；智能水雷通过大数据识别技术准确判别敌我舰船，实现自主攻防决策。
- **军民融合**：民企参与研发的北斗三号芯片，显著提升定位精度与数据处理能力，为系统提供强大底层算力支持。

<https://baijiahao.baidu.com/s?id=1841881729958096787&wfr=spider&for=pc>

<https://baijiahao.baidu.com/s?id=1841663536417320849&wfr=spider&for=pc>

目录

Background

Course Introduction

Interactive Data Exploration



学习目标



知识 目标

能够阐述清楚大数据、大数据分析实践的必要性，了解本门课程的考核方式及要求

重点

能力 目标

可根据课程要求规划本学期项目实施计划

重点

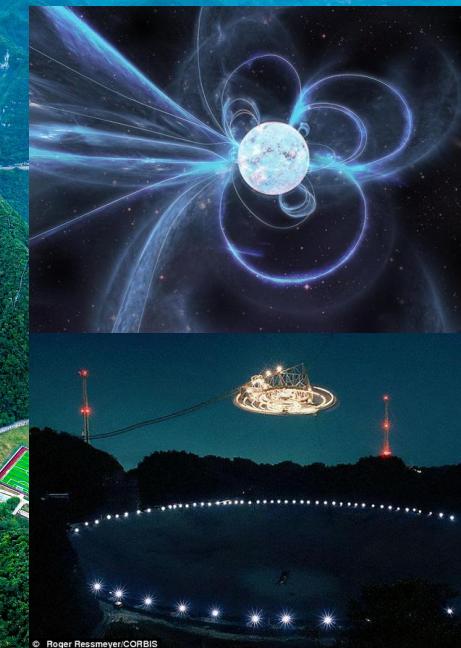
素质 目标

能够以批判性思维以及发展的眼光看待大数据分析

难点

世界最大：中国天眼FAST Five-hundred-meter Aperture Spherical radio Telescope

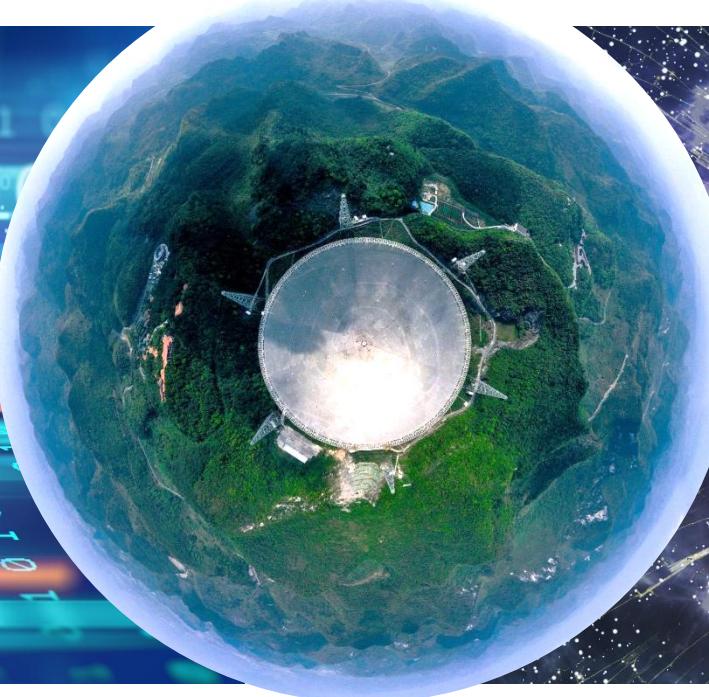
- 4450块反射面板单元
- 突破了射电望远镜的百米极限
- 拥有30个足球场大的接收面积
- 与号称“地面最大的机器”的德国波恩100米望远镜相比灵敏度提高约10倍





中国天眼FAST

数据量庞大
峰值数据每秒可以达到
38GB。
每年新增约10PB数据，
预计未来五年的数据总
量将超过100PB。



有效信息不明
捕捉到的海量原始数
据本身不能直接告诉
天文学家哪些是人类
未知的天文现象。

如何探索和处理海量数据以挖掘有效信息？

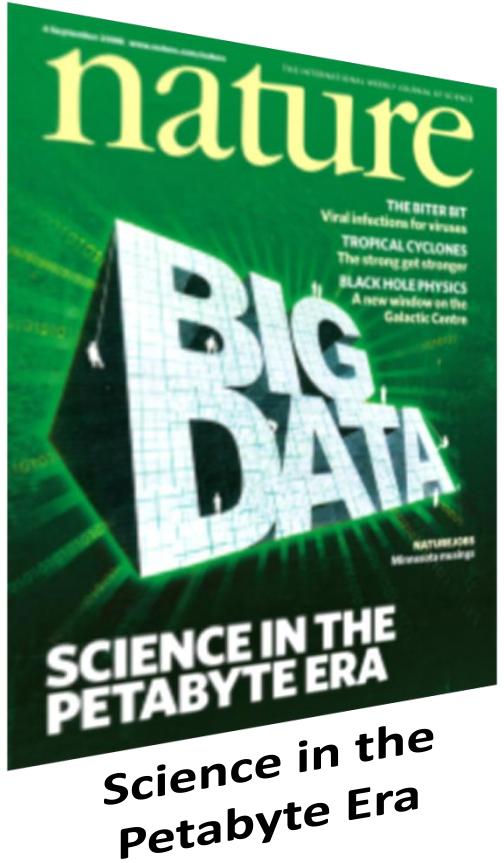


请联系阅兵式中可能有的大数据分析技术，
提2-5个关键词，描述什么数据是大数据？

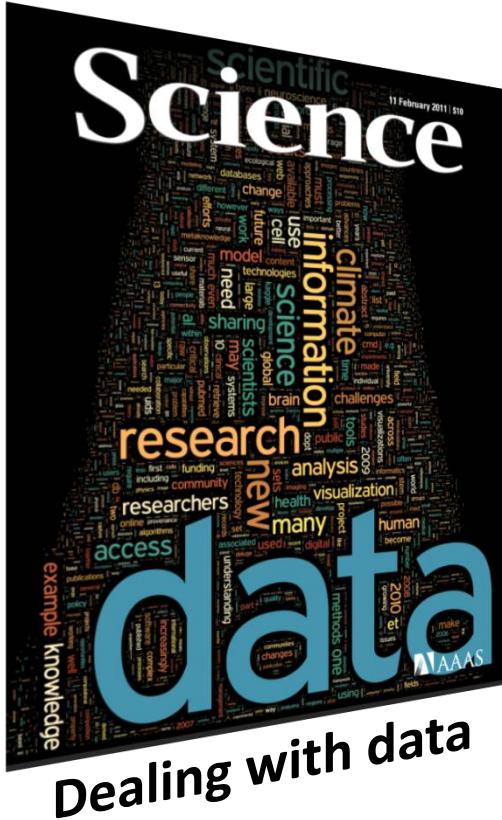
作答

Definition of Big Data

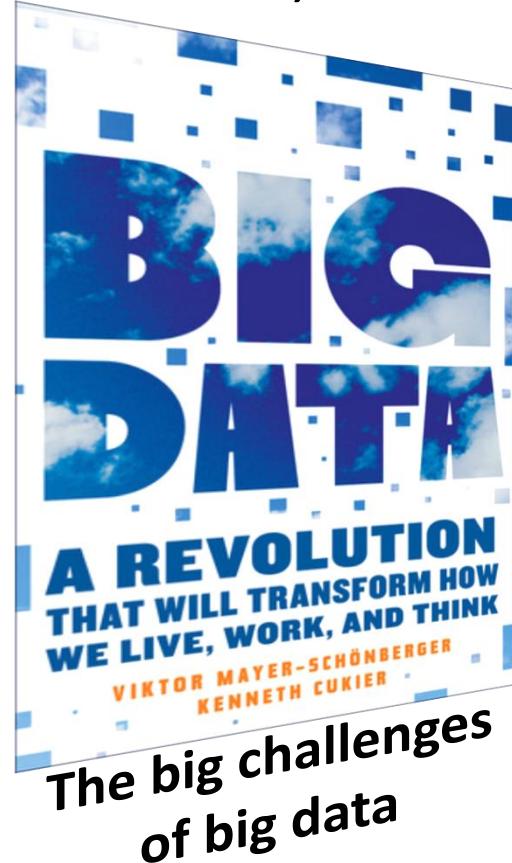
Nature 2008



Science 2011

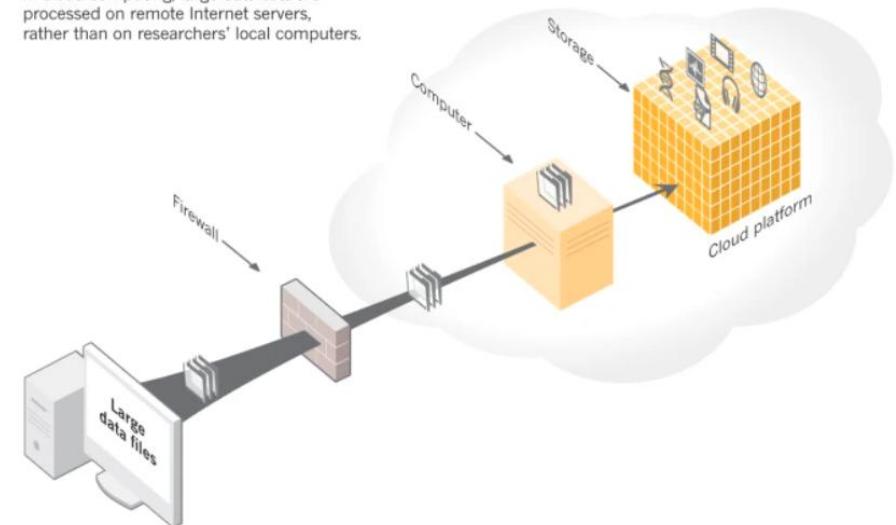


Nature, 2013



HEAD IN THE CLOUDS

In cloud computing, large data sets are processed on remote Internet servers, rather than on researchers' local computers.

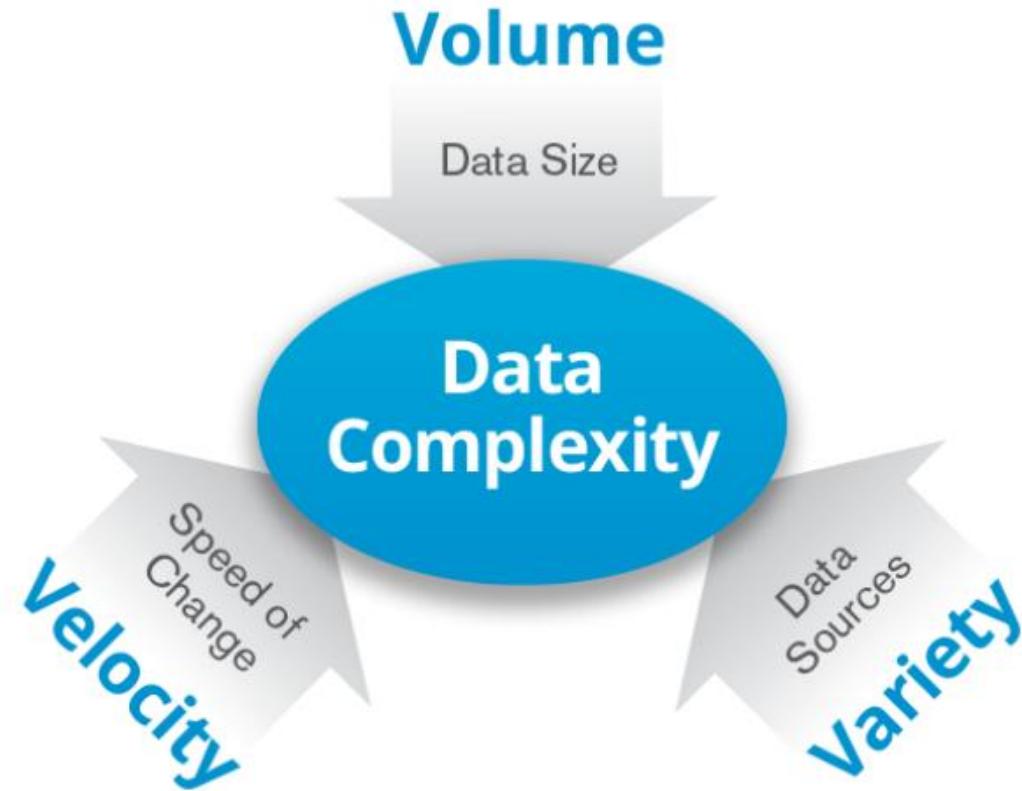




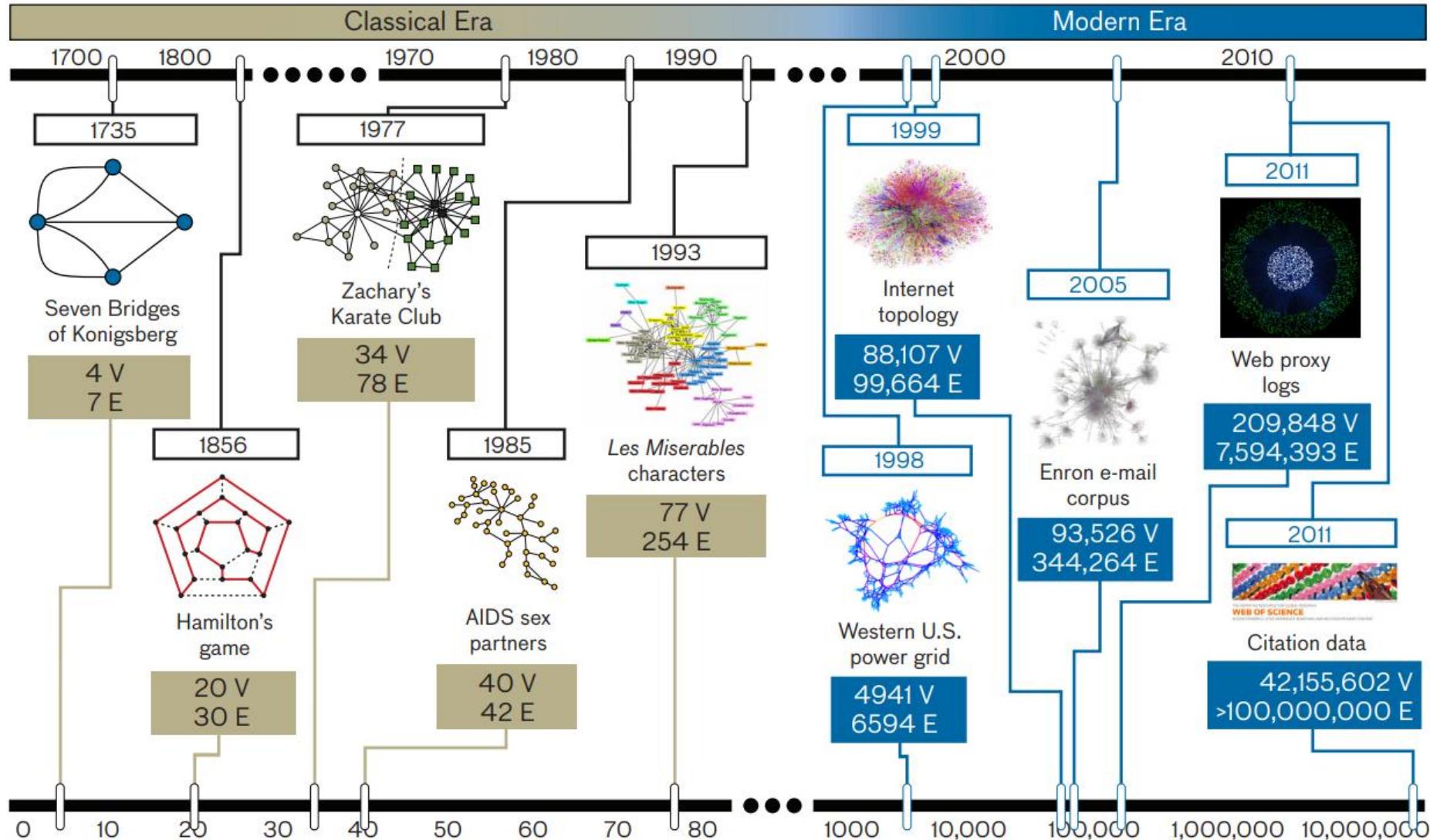
Definition of Big Data

*“Big data is **high-volume**, **high-velocity** and **high-variety** information assets that demand **cost-effective**, **innovative** forms of information processing for **enhanced insight** and **decision making**.”*

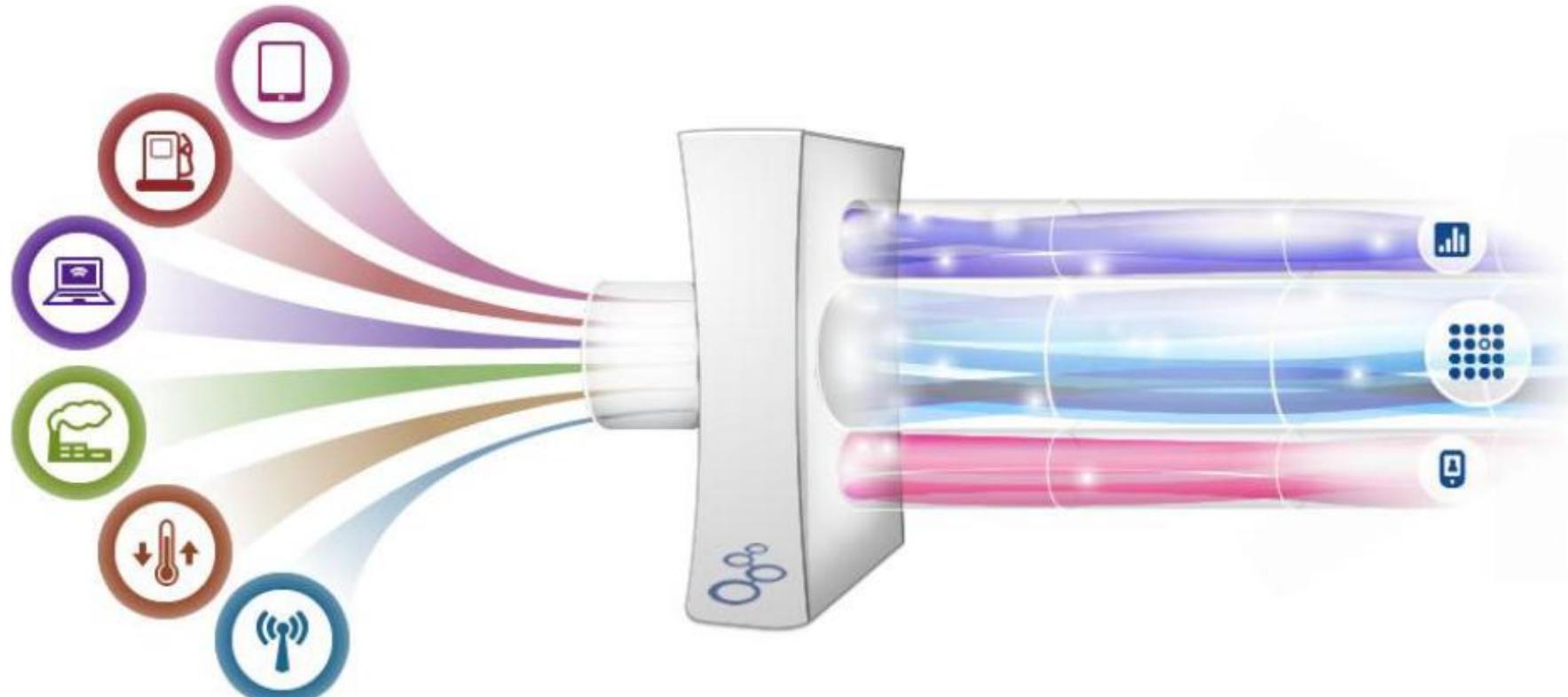
Gartner®



High Volume



High Velocity



① Real-time Events

② Continuous Analysis

③ Streaming Integration

High Variety



多源异构

海量高维

动态模糊

全局稀疏

局部冗余

时空断裂



Why Big Data?

- More data are being collected
- Open source code
- Commodity hardware / Cloud
- Domain requirements

- **High-Volume**
- **High-Velocity**
- **High-Variety**



Data Analytics

Improving human
intelligence with
machine intelligence

Big Data Analytics



IMDb: Movie Data

Find Movies, TV shows, Celebrities and more... All

Movies, TV & Showtimes Celebs, Events & Photos News & Community Watchlist

A movie poster for "It's a Wonderful Life" (1946). It features a man in a blue shirt and dark trousers carrying another man in a yellow shirt on his shoulders. The title "IT'S A WONDERFUL LIFE" is written in large red letters across the bottom. Other text on the poster includes "LIBERTY FILMS", "Frank CAPRA", "James STEWART", "Donna REED", and the names of other cast members like Lionel Barrymore, Thomas Mitchell, and Henry Travers.

It's a Wonderful Life (1946)  **Top 5000**

Approved 130 min - Drama | Family | Fantasy - 7 January 1947 (USA)

Your rating: ★★★★★★★★★★ 8.7 /10

Ratings: 8.7/10 from 202,743 users
Reviews: 632 user | 187 critic

An angel helps a compassionate but despairingly frustrated businessman by showing what life would have been like if he never existed.

Director: Frank Capra
Writers: Frances Goodrich (screenplay), Albert Hackett (screenplay), 4 more credits »
Stars: James Stewart, Donna Reed, Lionel Barrymore | See full cast and crew »

+ Watchlist Watch Trailer Share...

Details

Country: USA

Language: English

Release Date: 7 January 1947 (USA) See more »

Also Known As: The Greatest Gift See more »

Filming Locations: California, USA See more »

Box Office

Budget: \$3,180,000 (estimated)

Opening Weekend: £49,845 (UK) (19 December 2008)

Gross: £682,222 (UK) (24 December 2010)

See more »

Company Credits

Production Co: Liberty Films (II) See more »

Show detailed company contact information on IMDbPro »

Technical Specs

Runtime: 130 min | 118 min (DVD edition)

Sound Mix: Mono (RCA Sound System)

Color: Color (colorized) | Black and White

Aspect Ratio: 1.37 : 1

See full technical specs »

假如你是一位数据分析师，请根据这些电影票房数据，提出一个你感兴趣的问题。

The image shows a screenshot of the IMDb movie page for "It's a Wonderful Life" (1946). The page includes the movie poster, basic info (Approved, 130 min, Drama, Family, Fantasy), user rating (8.7/10 from 202,743 users), plot summary, director (Frank Capra), writers (Frances Goodrich, Albert Hackett), stars (James Stewart, Donna Reed, Lionel Barrymore), and links to the full cast and crew. On the right, there are sections for Details (Country: USA, Language: English, Release Date: 7 January 1947 (USA), etc.), Box Office (Budget: \$3,180,000, Opening Weekend: £49,845 (UK), etc.), Company Credits (Production Co: Liberty Films (II)), and Technical Specs (Runtime: 130 min | 118 min (DVD edition), Sound Mix: Mono (RCA Sound System), etc.).

作答

Big Data Analytics



IMDb: Movie Data



The image shows a screenshot of the IMDb movie page for "It's a Wonderful Life". The page features a large banner image of the movie poster, which includes the title "IT'S A WONDERFUL LIFE", the lead actors' names (James Stewart and Donna Reed), and the director's name (Frank Capra). Below the banner, there is a summary of the movie: "An angel helps a compassionate but despairingly frustrated businessman by showing what life would have been like if he never existed." The page also displays the movie's rating (8.7/10), the director (Frank Capra), writers (Frances Goodrich and Albert Hackett), stars (James Stewart, Donna Reed, Lionel Barrymore), and links to the full cast and crew. At the bottom, there are buttons for "Watchlist", "Watch Trailer", and "Share...".

Find Movies, TV shows, Celebrities and more... All

Movies, TV & Showtimes Celebs, Events & Photos News & Community Watchlist

It's a Wonderful Life (1946) Top 5000

Techniques to implement the analytics

An angel helps a compassionate but despairingly frustrated businessman by showing what life would have been like if he never existed.

Director: Frank Capra

Writers: Frances Goodrich (screenplay), Albert Hackett (screenplay), 4 more credits »

Stars: James Stewart, Donna Reed, Lionel Barrymore | See full cast and crew »

+ Watchlist Watch Trailer Share...

Details

Country: USA

Language: English

Release Date: 7 January 1947 (USA) See more »

Also Known As: The Greatest Gift See more »

Filming Locations: California, USA See more »

Box office

Budget: \$3,180,000 (estimated)

Opening Weekend: £49,845 (UK) (19 December 2008)

Gross: £68,111,411 (UK) (24 December 2010)

December 2008

Company Credits

Production Co: Liberty Films (II) See more »

Show detailed company contact information on IMDbPro »

Technical Specs

Runtime: 130 min | 118 min (DVD edition)

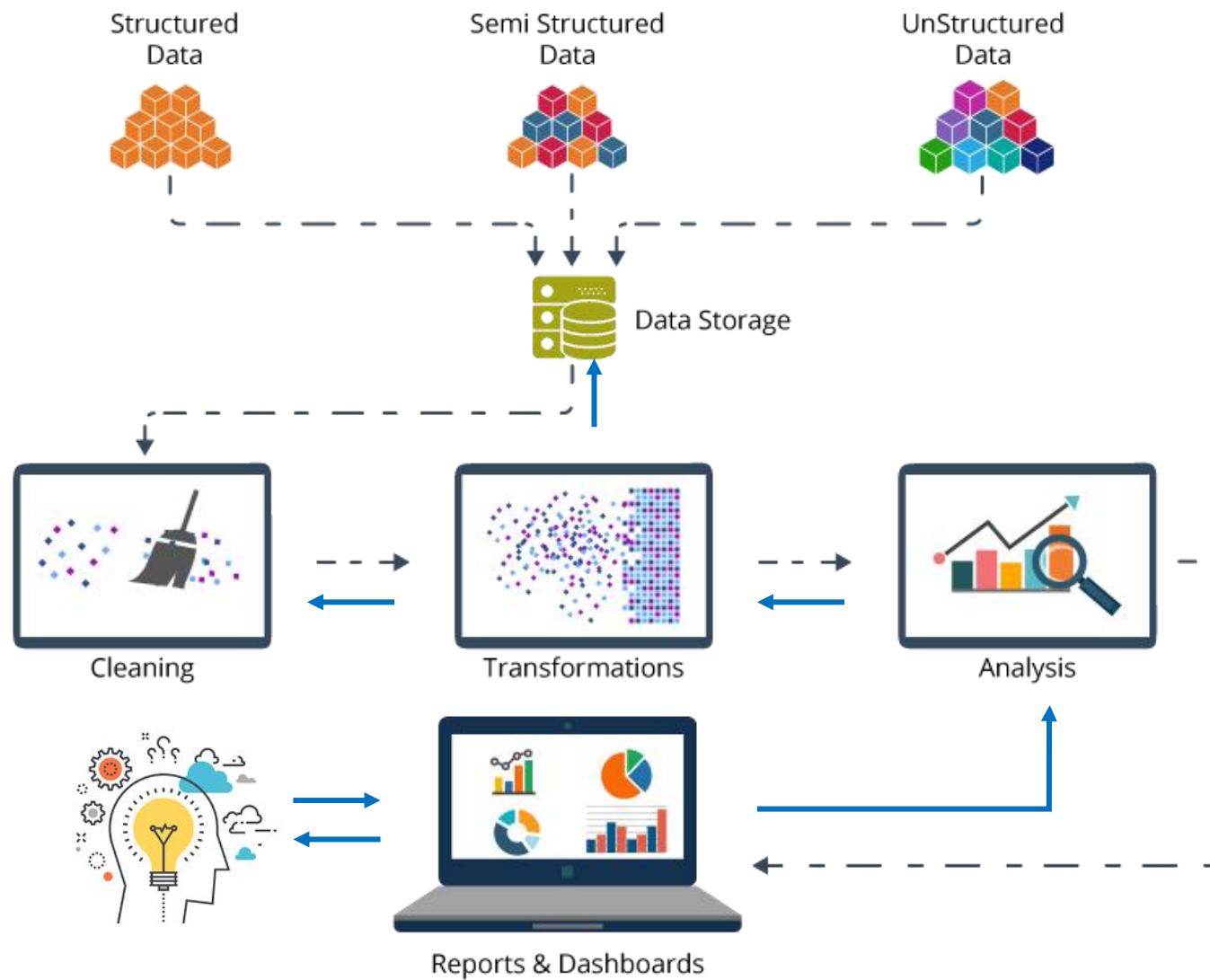
Sound Mix: Mono (RCA Sound System)

Color: Color (colorized) | Black and White

Aspect Ratio: 1.37 : 1

See full technical specs »

Big Data Analytics



- Huge Data Volumes Storage
- High-Speed Networks
- High Performance Computing
- Massive Parallelism
- Data Management
- Data Wrangling
- Data Mining and Analytics
- Machine Learning
- Data Visualization and interaction



Techniques Towards Big Data

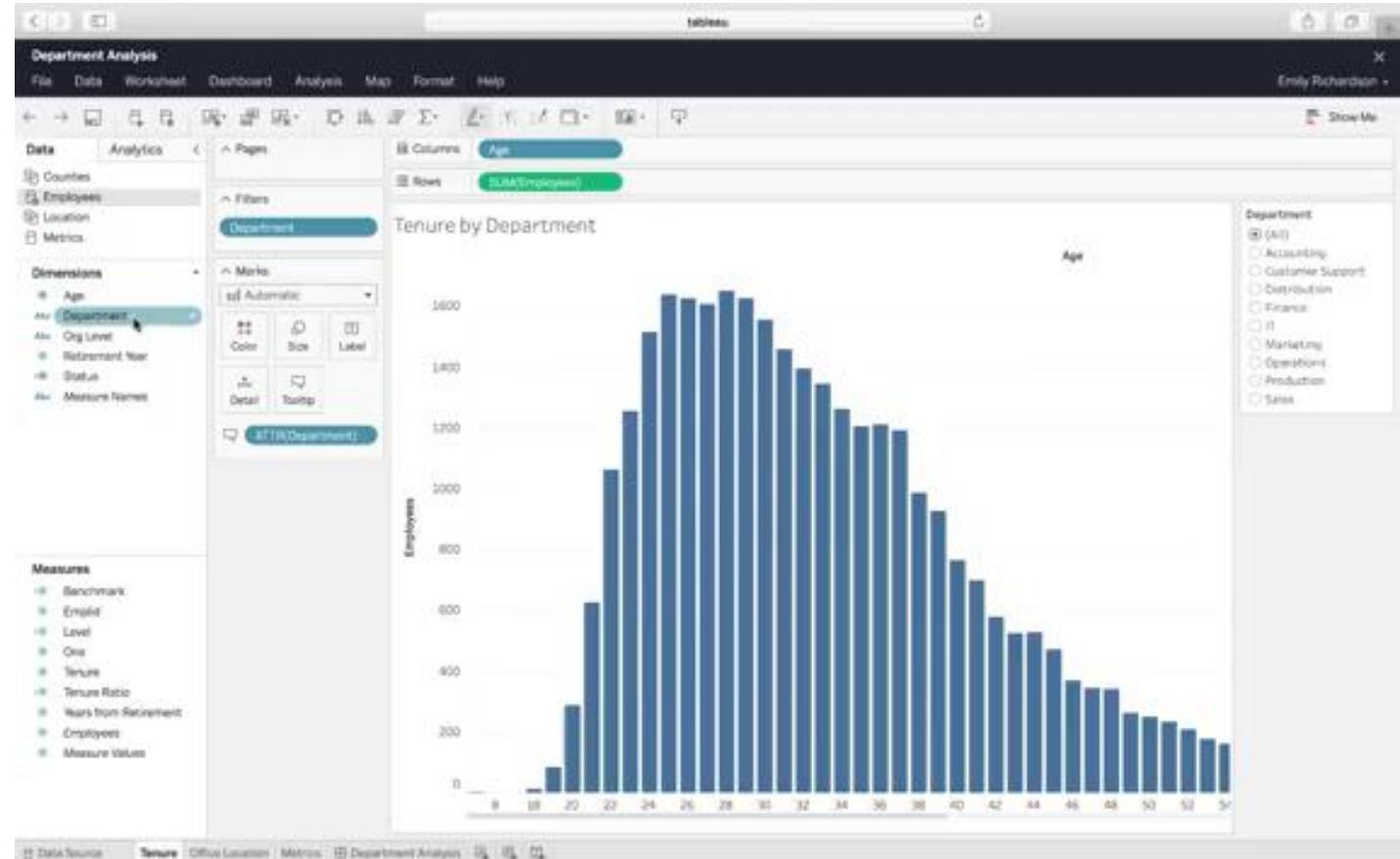
Interactive Data Exploration

- Huge Data Volumes Storage
- High-Speed Networks
- High Performance Computing
- Massive Parallelism
- Data Management
- Data Wrangling
- Data Mining and Analytics
- Machine Learning
- Data Visualization and interaction



Data Analytics

Visual Analytics



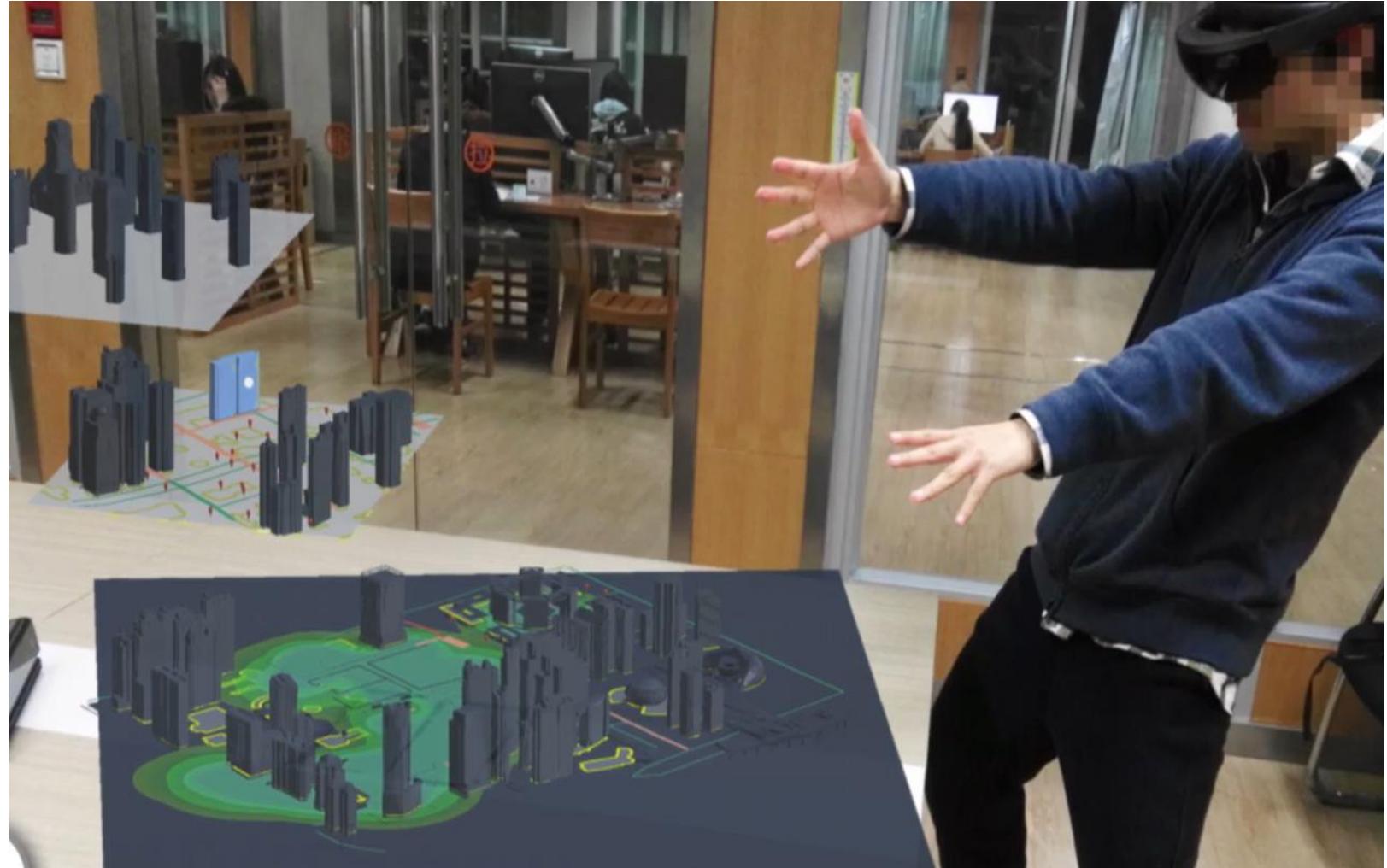
Harness the power of your data. **Unleash the potential of your people.**
Choose the analytics platform that disrupted the world of business intelligence.



Data Analytics



Immersive
Environment



Project-focused Course



课程背景

- 人机智能融合的大数据分析是国家安全情报决策、武器研究、智慧城市决策等国家战略领域及社会民生需求的基础技术手段
- 国家《“十四五”大数据产业发展规划》强调了数据分析对于构建稳定高效产业链的作用
- 然而我国大数据产业存在技术支撑不强，基础软硬件、开源框架与国际先进水平存在差距



史上最严技术出口管制

2018年11月，美国最新14类技术出口管制之第六条：数据分析（可视化、自动分析算法）；

(i) Systems-on-Chip (SoC); or
(ii) Stacked Memory on Chip.
(5) Advanced computing technology, such as:
(i) Memory-centric logic
(6) Data analytics technology, such as:
(i) Visualization;
(ii) Automated analysis algorithms; or
(iii) Context-aware computing.
(7) Quantum information and sensing technology, such as
(i) Quantum computing;
(ii) Quantum encryption; or
(iii) Quantum sensing.
(8) Logistics technology, such as:
(i) Mobile electric power;
(ii) Modeling and simulation;
(iii) Total asset visibility; or
(iv) Distribution-based Logistics

我国与美国之间的差距

大规模数据分析平台几乎被美国国家实验室所垄断

The slide compares the state of big data technology in China and the United States. It highlights the 'strictest export control' imposed by the US in November 2018, which specifically targets data analysis technologies like visualization and automated analysis algorithms. It also points out the significant gap in large-scale data analysis platforms, noting that they are almost entirely controlled by US national laboratories. Logos for various platforms like Oracle Database, SQL Server, and several visualization tools are shown.



课程历史

“人工智能+”引发**产业界**
大数据人才需求上涨

四名具有**国家级项目**经验
专业教师实施课程教学改革



- 2022 ● **大数据分析实践国际课程教改项目**
- 2021 ● 推出低代码声明反馈式**教学平台**
- 2020 ● 数据科学与大数据相关**省级教改项目**
- 2019 ● **重点教改**培育本课程，引入**国家级项目**作为课程实践项目
- 2018 ● 建设**数据科学与大数据技术**新工科创新试验班
- 2017 ● 教育部“新工科”建设
- 2016 ● 培养数据科学方向学生
- 2001 ● 《数据分析技术》《数据仓库数据挖掘》等大数据分析前身课程





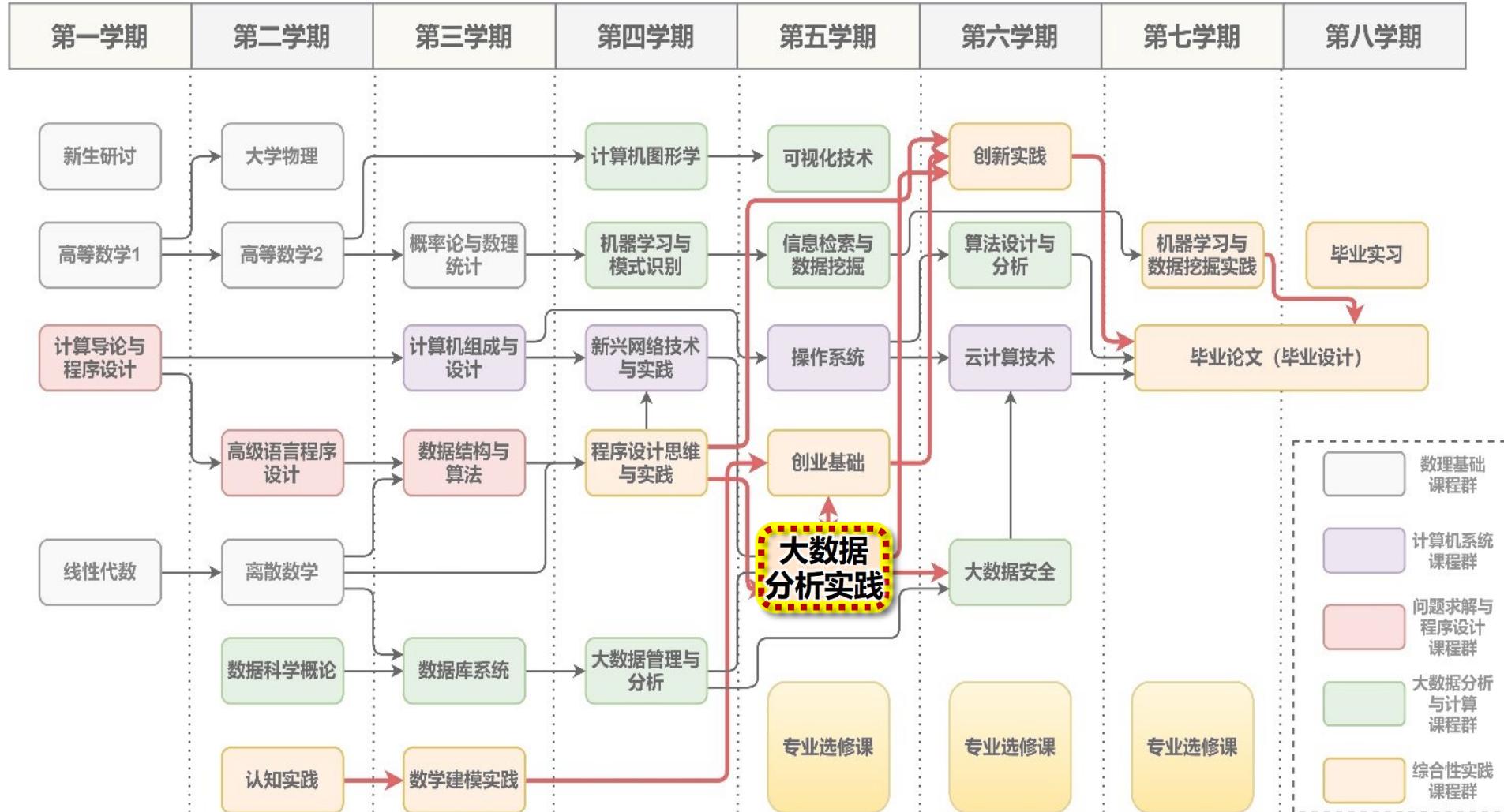
课程地位

专业培养目标

培养数据科学与大数据技术**拔尖**人才，学生具有坚实的数理基础和宽广扎实的计算机科学知识，具有**独立的研究能力**，熟练的**沟通能力**，具有逻辑推理、计算分析、算法优化、随机运用等方面的能力，具备理论思维、计算思维、数据思维、并行思维、实验思维等**科学素养**，实现**科学与工程**的紧密联系，具有高度的社会责任感和良好的职业道德，具有过硬的社会竞争力和国际化视野，具有不断**学习的能力和开拓创新**精神，具有良好的**团队合作和组织管理**能力，能够在业界工作中发挥领袖作用。



课程地位



课程类型

专业必修课

学时安排

课堂32学时
实验32学时

课程学分

3学分

面向专业

数据科学与大
数据技术专业

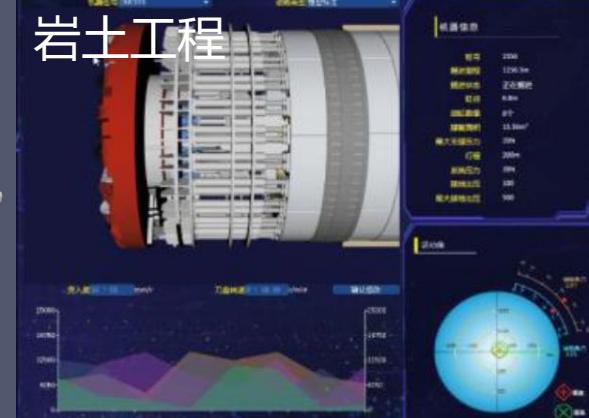
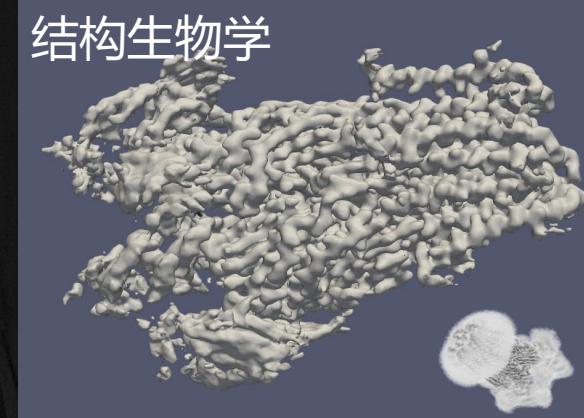
授课对象

本科三年级



课程职责

课程职责：培养学生构建人机智能融合的**大数据探索分析系统**能力





课程目标

课程职责：培养学生构建人机智能融合的**大数据探索分析系统**能力

素养
目标

- 树立为社会民生和国家战略研发国产大数据分析软件系统的**自主创新**意识
- 具备系统研发团队**合作精神**以及**职业素养**

能力
目标

- 自主研发面向**社会民生**以及**国家需求**的人机智能融合大数据分析系统
- 具备需求分析、前沿技术调研、系统开发、系统验证以及持续改进的**数据系统工程研究及实践能力**

知识
目标

- 掌握数据管理、数据分析、数据可视化、人机交互等开发大数据探索式分析系统所需的**算法设计思想**、**基础方法**、**核心技术及实现流程**

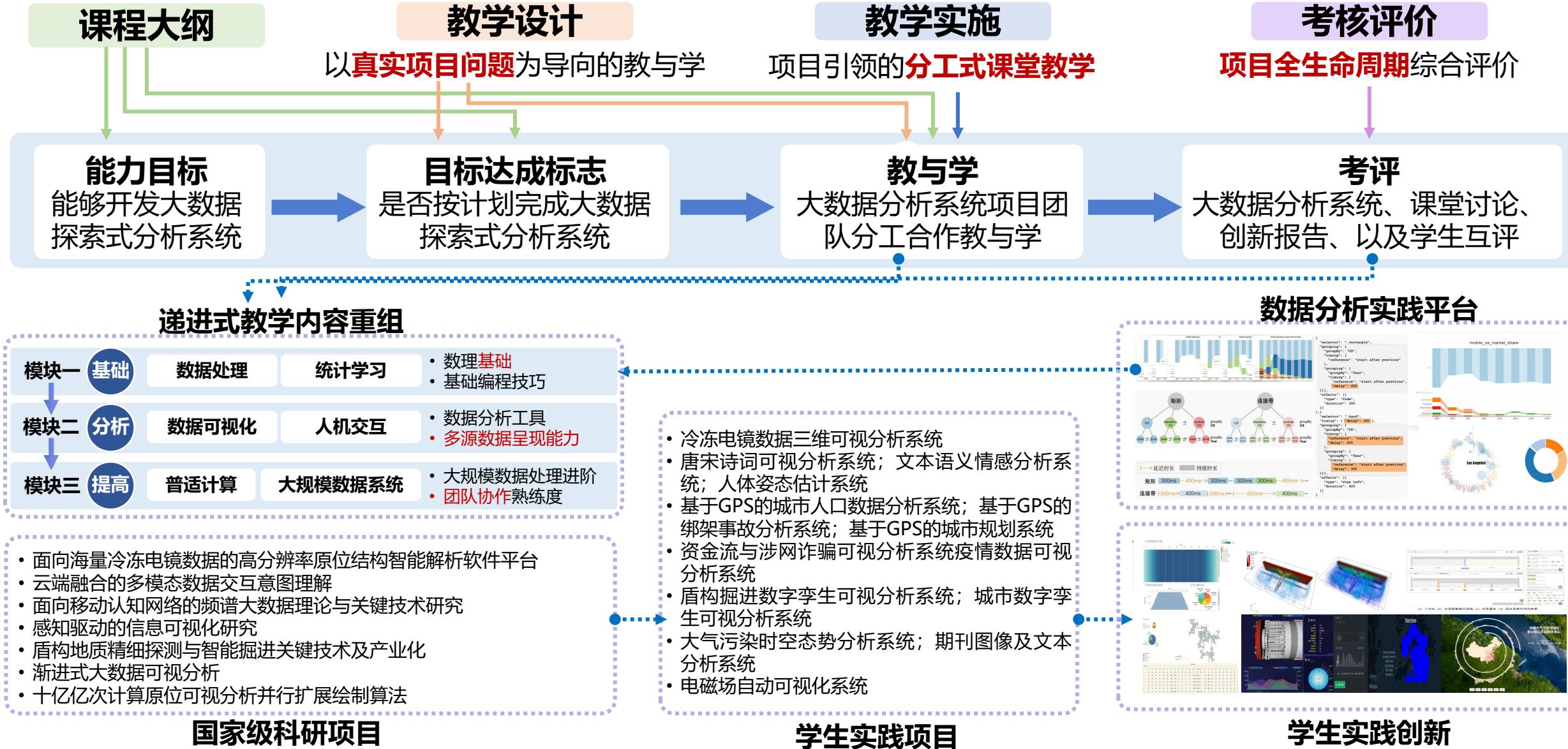
- 关注社会民生与国家战略
- 自主创新
- 团结合作

- 系统实践能力
- 系统研究能力

- | | |
|-----------|---------|
| • 数据处理 | • 数据可视化 |
| • 统计学习 | • 人机交互 |
| • 大规模数据系统 | • 普适计算 |



目标细化





课程内容

目标：构建大数据探索式分析系统



社会民生和国家战略为载体的 真实项目解构

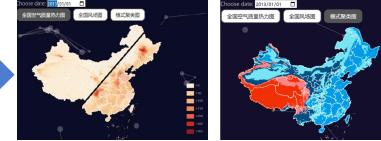
项目来源	国家项目名称	实践项目名称
国家重点研发计划	面向海量冷冻电镜数据的高分辨率原位结构智能解析软件平台	冷冻电镜数据三维可视分析系统
国家重点研发计划	云端融合的多模态数据交互意图理解	唐宋诗词可视分析系统 文本语义情感分析系统 人体姿态估计系统
国家重点研发计划	面向移动认知网络的频谱大数据理论与关键技术研究	基于GPS的城市人口数据分析系统 基于GPS的绑架事故分析系统 基于GPS的城市规划系统
国家自然科学基金	感知驱动的信息可视化研究	资金流与涉网诈骗可视分析系统 疫情数据可视分析系统
山东省重点研发计划项目	盾构地质精细探测与智能掘进关键技术及产业化	盾构掘进数字孪生可视分析系统 城市数字孪生可视分析系统
科技部高端外国专家引进计划	渐进式大数据可视分析	大气污染时空态势分析系统 期刊图像及文本分析系统
国防科工局国防基础科研核技术专项	原位可视分析并行扩展绘制算法 十亿级次计算	电磁场自动可视化系统

大气污染实时态势分析系统

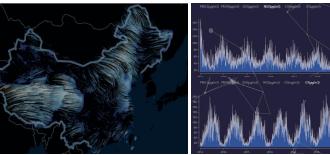
基于真实项目执行过程的理论教学知识重构

理论配套小实验

探索交互分析



呈现数据



建立数据模型

LSTM + K-Means

Initial preparation
• As the daily fluctuation greatly, today is selected as the window size.
• Training set: test set: 8:2
• Data processing: separate the month and day and convert them to int format
`data_ae_baseline`

	AQI	PM2_5	PM10	SO2	NO2	CO	O3	Time	Year	Month	Day
0	75.25000	55.20000	65.41004	27.07	25.50001	1.11	27.40669	2013-01-01	2013	1	1
1	117.22498	88.77996	101.50000	30.91	32.83002	1.61	23.76001	2013-01-02	2013	1	2
2	86.90002	45.20000	55.77000	22.71	24.00000	1.81	26.12986	2013-01-03	2013	1	3
3	78.37500	57.70001	55.86001	20.68	21.32000	1.56	26.02000	2013-01-04	2013	1	4
4	125.60000	95.48003	81.22001	27.82	28.28001	1.42	19.90000	2013-01-05	2013	1	5

收集及处理数据

中国大气污染源数据集

- 时间: 2013-2018
- 空间分辨率: 15 kilometers
- 空间坐标: 经度、纬度
- 污染源变量: PM2.5, PM10, O₃, CO, SO₂, NO₂
- 大气变量: 风向U、风向V、温度、相对湿度、压强
- 40000行数据、13个维度

理解项目需求

分析任务:

- 大气污染源分析
- **大气污染时空态势分析**
- 大气污染传输模式分析
- 大气污染预测

普适计算

提高

大规模计算

可视化设计

可视化评估及实践工具

人机交互、Canis平台

分析

统计分析

机器学习与深度学习工具

数据采样与降维

数据质量管理

众包

电子表格

基础

大数据分析系统

实践项目导论

科研实践入门

系统构建经验谈

SPARK实践

手机移动数据采集与分析

人机交互实践

BERT实践

BERT环境配置

可视化设计

统计方法实践

电子表格实践

基于Pandas的数据清洗

基于Pandas的数据采样

项目规划

团队分组

课程大纲

6 Personal assignments

6 Group assignments



上课日期	授课内容	实验内容	周次
20250903	课程入门、大数据探索式分析	课程实践项目介绍、项目组队测试	第一周
20250910	项目经验谈、科研实践入门	项目管理工具制定项目计划	第二周
20250917	数据采样与降维	数据采样实践	第三周
20250924	数据质量管理	数据质量实践	第四周
20251001	众包与电子表格	电子表格实践	第五周
20251008	统计分析方法与工具	统计方法实践	第六周
20251015	可视化设计	可视化设计实践	第七周
20251022	中期汇报（论文+项目进展）1	中期进展报告	第八周
20251029	中期汇报（论文+项目进展）2	BERT实践环境配置	第九周
20251105	机器学习方法与工具	BERT实践	第十周
20251112	人机交互方法与工具	Canis/Cast/Libra实践	第十一周
20251119	普适计算	手机移动数据采集与分析	第十二周
20251126	大规模数据分析系统	SPARK实践	第十三周
20251202	如何撰写项目论文	大项目收尾	第十四周
20251209	项目结题报告1		第十五周
20251216	项目结题报告2	大项目验收	第十六周



Basic Factors



No exams!

Faculty Qiong Zeng (曾琼)

Teaching Assistant Xi Duan (段曦)

Zhiyuan Meng (孟致远)

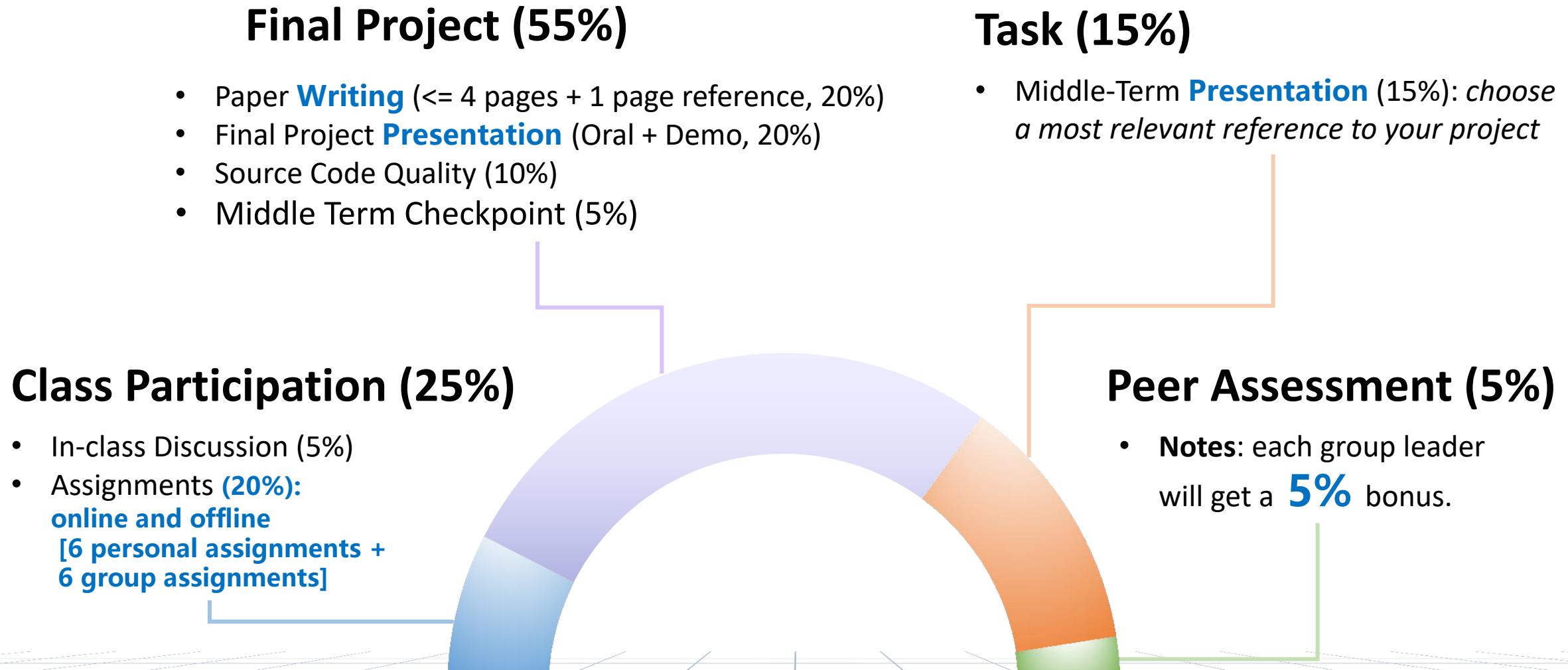
Time Wednesday 10:10-12:00 (振声苑N202)

Friday 16:10-18:00 (N3)

Office Hour Friday 10:00-12:00 (N3-412)



Course Grading





Course Grading

一级指标		二级指标		优秀	良好	中等	及格	不及格
指标	比例	指标	比例	≥9分	8-9分	7-8分	6-7分	<6分
大数据分析系统演示	20%	系统完整性	40%	系统功能完整，涵盖数据清洗、数据分析挖掘、数据可视化、数据交互等模块，能够满足任务分析的所有需求。	系统功能基本完整，基本涵盖数据清洗、数据分析挖掘、数据可视化、数据交互等模块，能够满足任务分析的大部分需求。	系统功能不够完整，涵盖数据清洗、数据分析挖掘、数据可视化、数据交互等模块有较大缺失，但能够满足任务分析的大部分需求。	系统功能不完整，涵盖数据清洗、数据分析挖掘、数据可视化、数据交互等模块有较大缺失，无法满足任务分析的多数需求，存在一些系统问题。	系统功能严重不完整，未涵盖数据清洗、数据分析挖掘、数据可视化、数据交互等模块，无法满足任务分析的基本需求。
		系统创新性	20%	系统采用了新的技术或方法，具有较强的创新性，并取得了显著的效果。	观点基本明确，论述基本清晰，基本易于理解；术语使用基本规范，解释基本到位；概念定义基本准确，无重大歧义；句子结构基本合理，语句通顺。	观点不够明确，论述部分清晰；术语使用不够规范，解释存在较大不到位；概念定义不够准确，有较多歧义；句子结构不够合理，阅读困难。	观点模糊，论述基本清晰，较难理解；术语使用基本规范；概念定义不准确，存在重大歧义；句子结构表述混乱。	观点不明确，论述无逻辑，术语使用错误，缺少解释，语句不通顺



Course Grading

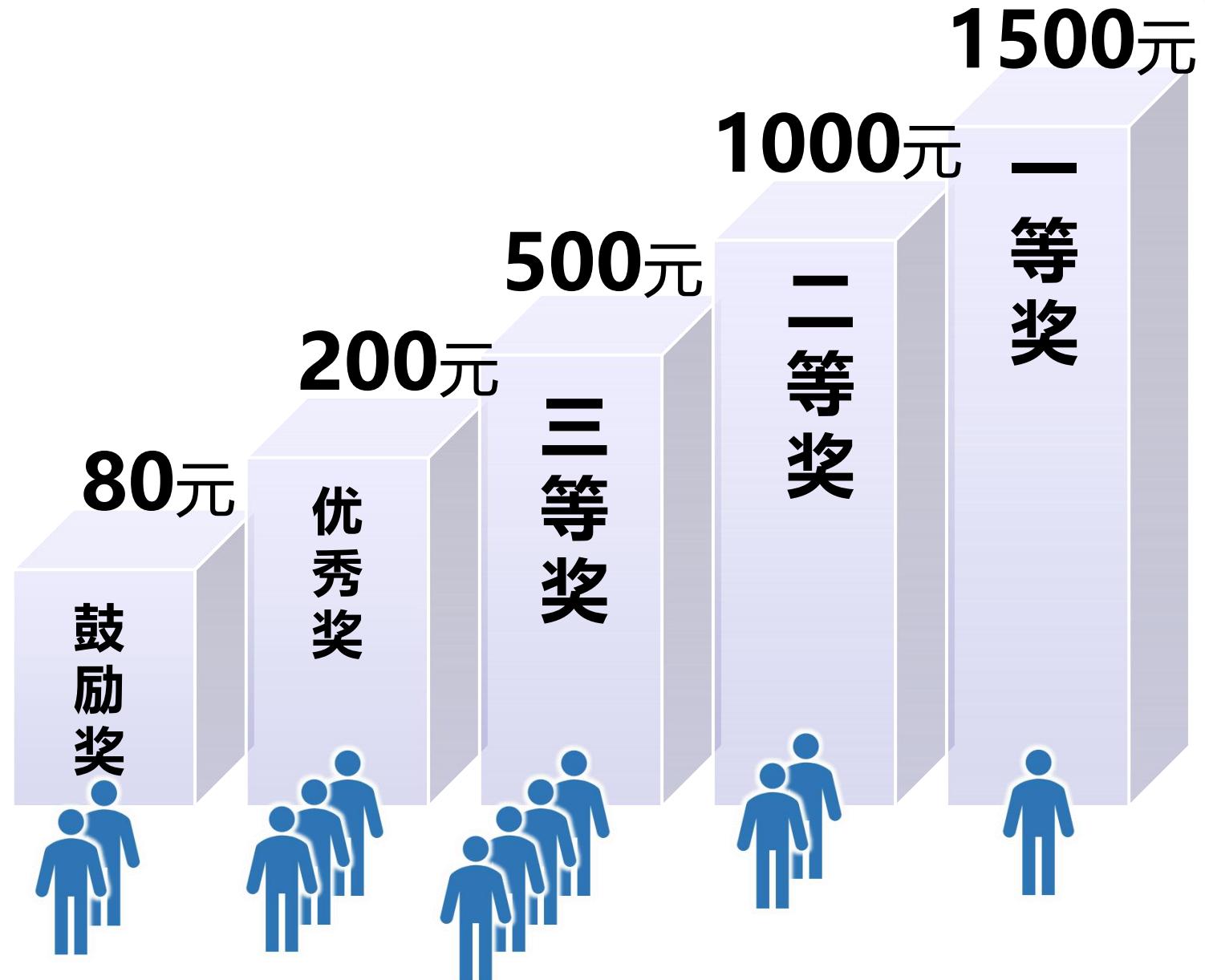
一级指标		二级指标		优秀	良好	中等	及格	不及格
指标	比例	指标	比例	≥9分	8-9分	7-8分	6-7分	<6分
项目论文	20%	写作规范	15%	论文结构完整，包含引言、相关工作、技术细节、实验分析、总结、参考文献等核心模块，各模块内容充实完整。语言表达准确、流畅，无语法错误；图表清晰美观，图标题描述自洽。	论文结构论文结构完整，包含引言、相关工作、技术细节、实验分析、总结、参考文献等模块，各模块均有涉及。语言表达基本准确，少量语法错误；图表清晰美观。	论文结构论文结构基本完整，参考文献、结论等模块存在缺失。语言表达基本准确，语法错误较多。图表描述不符合规范。	论文结构论文结构缺失参考文献部分。语言表达基本准确，语法错误较多。图表描述不清晰，排版存在一定混乱。	论文结构存在严重缺失，比如无相关文献调研。语言表达存在较大逻辑错误以及逻辑不通顺。无图表或图表混乱。
		语言表达	20%	观点明确，论述清晰，易于理解；术语使用规范，解释清楚到位；概念定义准确，无歧义；句子结构合理，语句通顺。	观点基本明确，论述基本清晰，基本易于理解；术语使用基本规范，解释基本到位；概念定义基本准确，无重大歧义；句子结构基本合理，语句通顺。	观点不够明确，论述部分清晰；术语使用不够规范，解释存在较大不到位；概念定义不够准确，有较多歧义；句子结构不够合理，阅读困难。	观点模糊，论述基本清晰，较难理解；术语使用基本规范；概念定义不准确，存在重大歧义；句子结构表述混乱。	观点不明确，论述无逻辑，术语使用错误，缺少解释，语句不通顺



Course Grading

一级指标		二级指标		优秀	良好	中等	及格	不及格
指标	比例	指标	比例	≥9分	8-9分	7-8分	6-7分	<6分
项目论文	20%	写作规范	15%	论文结构完整，包含引言、相关工作、技术细节、实验分析、总结、参考文献等核心模块，各模块内容充实完整。语言表达准确、流畅，无语法错误；图表清晰美观，图标题描述自洽。	论文结构论文结构完整，包含引言、相关工作、技术细节、实验分析、总结、参考文献等模块，各模块均有涉及。语言表达基本准确，少量语法错误；图表清晰美观。	论文结构论文结构基本完整，参考文献、结论等模块存在缺失。语言表达基本准确，语法错误较多。图表描述不符合规范。	论文结构论文结构缺失参考文献部分。语言表达基本准确，语法错误较多。图表描述不清晰，排版存在一定混乱。	论文结构存在严重缺失，比如无相关文献调研。语言表达存在较大逻辑错误以及逻辑不通顺。无图表或图表混乱。
		语言表达	20%	观点明确，论述清晰，易于理解；术语使用规范，解释清楚到位；概念定义准确，无歧义；句子结构合理，语句通顺。	观点基本明确，论述基本清晰，基本易于理解；术语使用基本规范，解释基本到位；概念定义基本准确，无重大歧义；句子结构基本合理，语句通顺。	观点不够明确，论述部分清晰；术语使用不够规范，解释存在较大不到位；概念定义不够准确，有较多歧义；句子结构不够合理，阅读困难。	观点模糊，论述基本清晰，较难理解；术语使用基本规范；概念定义不准确，存在重大歧义；句子结构表述混乱。	观点不明确，论述无逻辑，术语使用错误，缺少解释，语句不通顺

Course BONUS





Expectations from You

- Be comfortable **asking questions, identifying a problem**
- Be comfortable **reading research papers**
- Be comfortable **coding**: data processing, machine learning, visualizations, interaction
- **Active participate** in discussions and our class
- Be comfortable summarizing your **progress**
- Be comfortable **writing up** experimental results in a **research paper format**

Project Topics



NEW DESIGN

	Topics	潜在产出	合作教师
1	IEEE SciVis Contest 海洋大气气候数据可视分析系统【冠军奖：1000刀+IEEE CG&A论文】		曾琼、段曦
2	基于华为昇腾算力与计图框架的冷冻电镜数据分析大模型构建与效能优化【华为云】	系统、论文、专利或软著	曾琼、王飞宇
3	融合知识图谱与时序分析的海洋新闻态势感知系统【国家深海基地合作】	系统、论文、专利或软著	曾琼、孟致远
4	基于实时位置数据的传染病密接个体识别系统【国自然项目】	系统、论文、专利或软著	滕德军、曾琼
5	大规模空间数据交叉匹配数据处理系统【国自然项目】		滕德军、曾琼
6	人类肠道微生物数据资源库与多维智能分析平台【微生物国重合作】	系统、专利、软著、	曾琼、蒋荷
7	《大数据分析实践》智慧课程平台【国家一流课程建设经费资助】	系统、专利、软著、	曾琼、栾峻峰



Project Groups

竞标: 队长 → 队员 → 项目

基本规则: 每组获得相同数量的初始金币 (10) , 金币将用于**队员的选取以及项目题目**



- 3题, 满分100, 30min
- 得分=平均百分制绩点*0.5+随堂测试得分*0.5
- 每个人的金币数量为 $0.1 * (\text{选课总人数} + 1 - \text{得分排名})$

- 队员总身价**不可超过**队伍初始金币
- **同时**选项目
- 出现多个队伍选择同样项目时, 根据团队**总身价先后**顺序排名
- 每个项目**少于等于2个**队伍选择

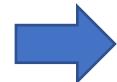


Project Groups

竞标：队长 → 队员 → 项目

基本规则：每组获得相同数量的初始金币（10），金币将用于队员的选取以及项目题目

随堂测试
选队长



双方互联

项目选取

- 3题，满分100，30min
- 得分=上学期百分制绩点*0.5+随堂测试得分*0.5
- 每个人的金币数量为 $0.1 * (\text{选课总人数} + 1 - \text{得分排名})$

周五下午赠
提神版随堂测试

！！！

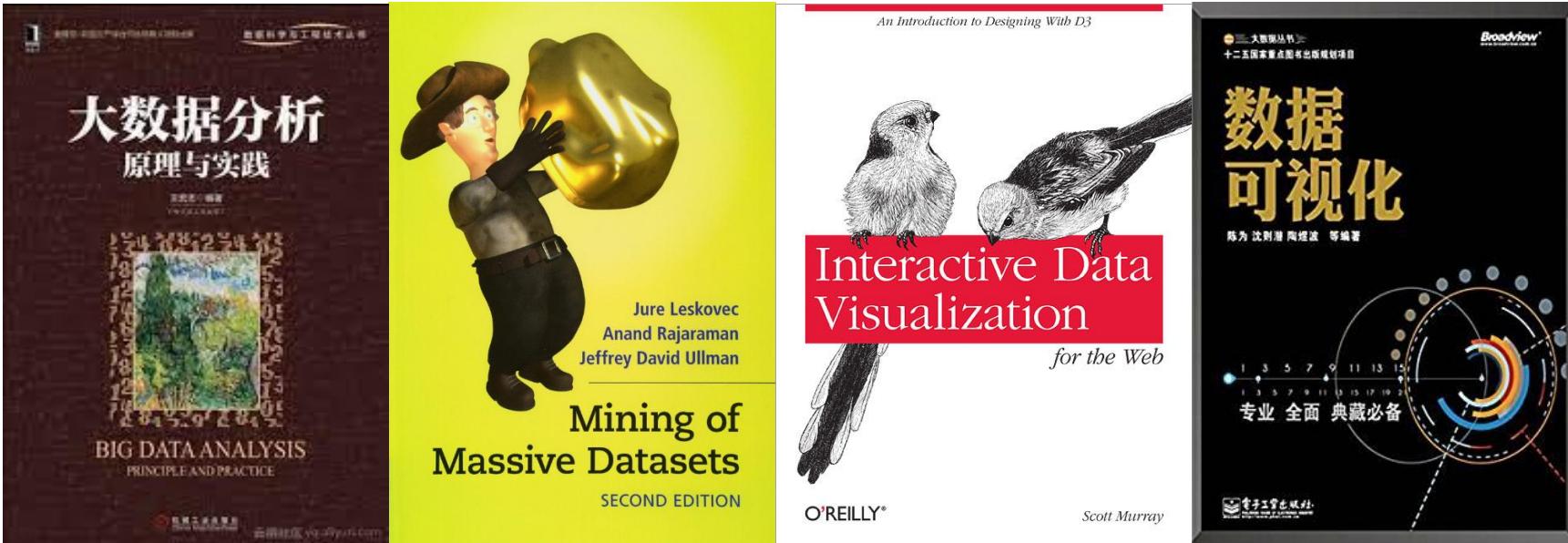
- 队员总身价不可超过队伍初始金币数
- 同时选项目优先级：各队伍项目选择先后顺序排名
- 总身价先按项目选择先后顺序排名
- 每个项目不超过3个队伍选择

Week																
	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	
G1	ZQ	DX	DX	ZY	ZY	ZY	DX	ZY	ZQ	DX	DX	ZY	ZY	ZQ	DX	
G2	ZQ	DX	DX	ZY	ZY	ZY	DX	ZY	ZQ	DX	DX	ZY	ZY	ZQ	DX	
G3	DX	ZQ	DX	ZY	ZY	ZY	DX	ZQ	DX	ZQ	DX	ZY	DX	ZY	ZY	
G4	DX	ZQ	DX	ZY	ZY	ZY	DX	ZQ	DX	ZQ	DX	ZY	DX	ZY	ZY	
G5	DX	ZY	ZQ	DX	DX	DX	ZQ	ZY	DX	ZY	ZQ	DX	ZY	ZY	ZY	
G6	DX	ZY	ZQ	DX	DX	DX	ZQ	ZY	DX	ZY	ZQ	DX	ZY	ZY	ZY	
G7	ZY	ZY	ZY	ZQ	DX	ZQ	ZY	DX	ZY	ZY	ZY	ZQ	DX	DX	DX	
G8	ZY	ZY	ZY	ZQ	DX	ZQ	ZY	DX	ZY	ZY	ZY	ZQ	DX	DX	DX	
G9	ZY	DX	ZY	DX	ZQ	DX	ZY	DX	ZY	DX	ZY	ZQ	ZY	ZQ		
G10	ZY	DX	ZY	DX	ZQ	DX	ZY	DX	ZY	DX	ZY	DX	ZQ	ZY	ZQ	
G11	ZY	ZY	ZY	ZQ	DX	ZQ	ZY	DX	ZY	CC	ZY	ZQ	DX	DX	DX	
G12	ZY	DX	ZY	DX	ZQ	DX	ZY	DX	ZY	DX	ZY	DX	ZQ	ZY	ZQ	



Course Information

课程网站: <https://sdubigdatacourse.github.io/#/>



- **Conference:** SIGMOD, KDD, VIS, CHI, SIGIR, MM,
- **Researchers (acknowledgements):**
Daniel Keim: *Mastering the information age*
Stratos Idreos: *Overview of Data Exploration Techniques*

学生成果



2019-2020-1	2020-2021-1	2021-2022-1	2022-2023-1	2023-2024-1
AR空间多类散点图采样与显示 语义情感分析系统 多人姿态估计与跟踪系统	面向珍惜鸟类保护的环境监测系统 工厂安全监测与分析 基于GPS绑架事故分析 公园内游客行为模式及异常行为分析	工业环境监测可视分析 城市监视器数据的公民救援分析 保护区鸟类模式分析 水道污染环境监测系统 大气污染可视分析系统	城市数据聚类分析 城市数字孪生可视分析 GIS数据城市绑架事故分析 城市数据安全监测系统 时空数据可视分析系统	大规模冷冻电镜数据可视分析系统 海洋非法捕捞识别系统 大气污染时空经济效益可视分析系统 全球温度趋势分析系统 计算机科学家社交网络分析 中医药理可视分析系统 全场景生鲜超市库存决策系统 运营商大数据栅格时序图预测系统
多类蓝噪声采样、AR/VR技术、LSTM、RAM、DenseNet	Tableau、DBSCAN聚类算法、参数化聚类算法、数据库、用户交互式可视化界面、CNN、五折交叉验证	D3、Tableau、无监督聚类算法、ARIMA模型、Dynamic Time Warping、LSTM	Tableau、K均值、DBSCAN、TOPSIS、word2Vec、GPT-2、LSTM、WMD、ArcGIS、PageRank、高斯和计算密度估计、KM算法、协同过滤算法、Vue3、Echarts、D3	线性回归、XGBoost、LSTM、D3、LightGBM、贪心算法、SVM、GCN、BiLSTM、逆地理编码、t-SNE、孤立森林、Streamlit、Transformer、随机森林方法、ECharts、GAN、Redis、Nginx、Docker、Vue.js、Three.js、LLAMA大模型、GPT-4、ChatGPT、熵权法、KNN、SIREN、ONNX.js



学生成果



大气污染时空经济效益可视分析系统
项目来源：生物气溶胶多组分高灵敏度在线监测技术（科技部）

后端：

编程语言：Python
Web 框架：FastAPI
数据库 ORM： SQLAlchemy
数据库： MySQL
GPT-4

测试：

接口测试工具：Hoppscotch

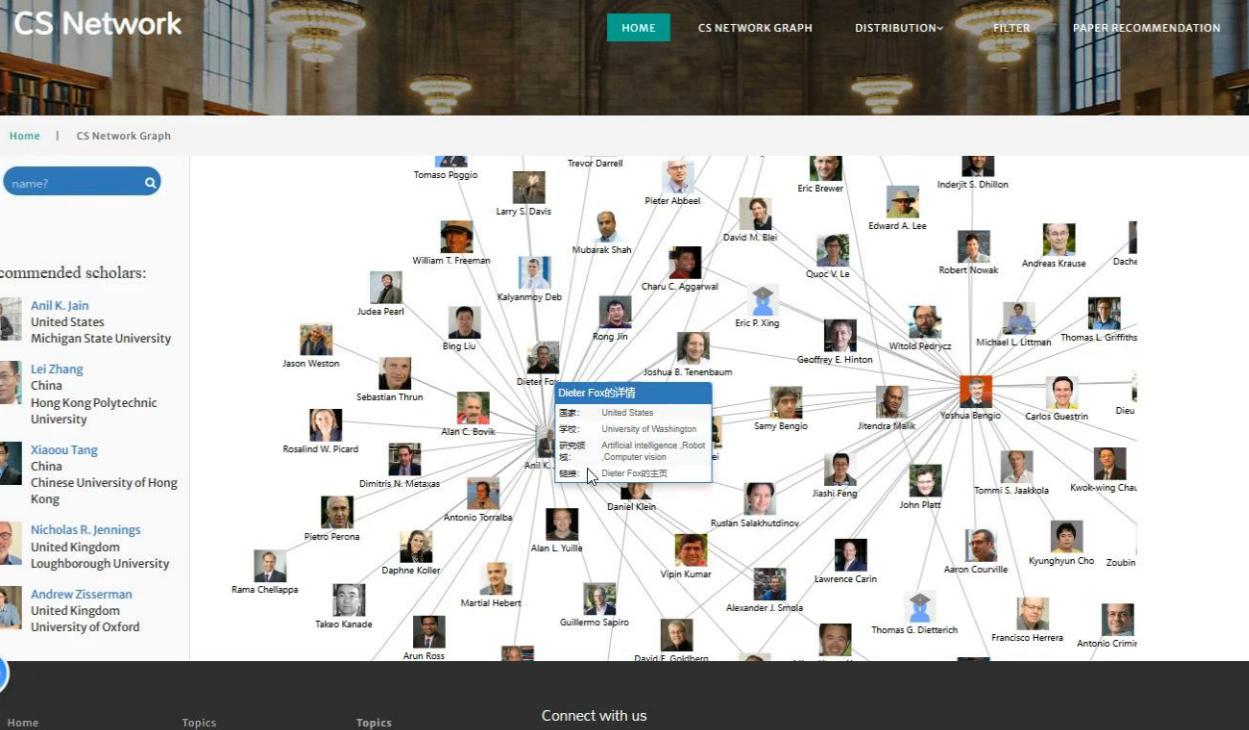
前端：

JavaScript 框架：Vue.js
可视化库：ECharts
地图 API：百度地图 API

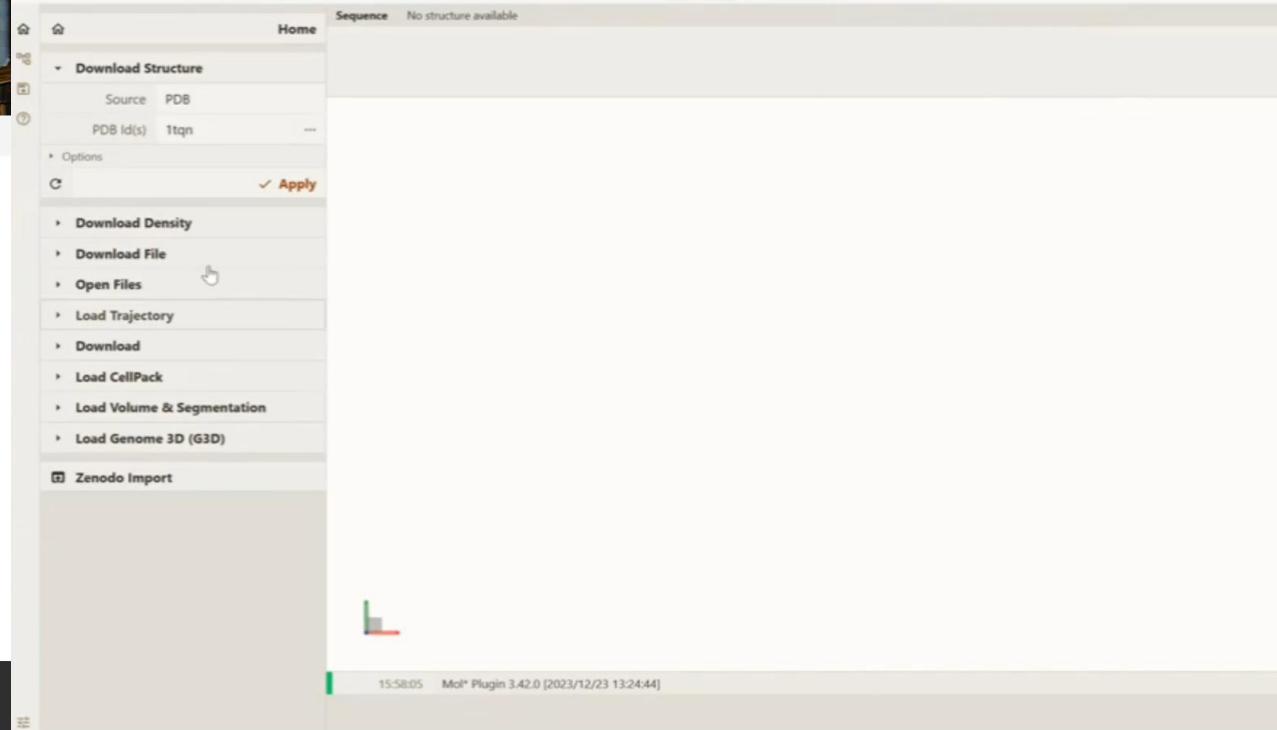
部署：

主从式架构
分布式文件系统：Minio
缓存与消息队列：Redis
负载均衡与反向代理：Nginx
容器虚拟化技术：Docker

学生成果



计算机科学社交网络分析
来源：渐进式大数据可视分析（科技部）



大规模冷冻电镜数据分析
来源：面向海量冷冻电镜数据的
高分辨率原位结构智能解析平台（科技部）



Practices on Big Data Analytics

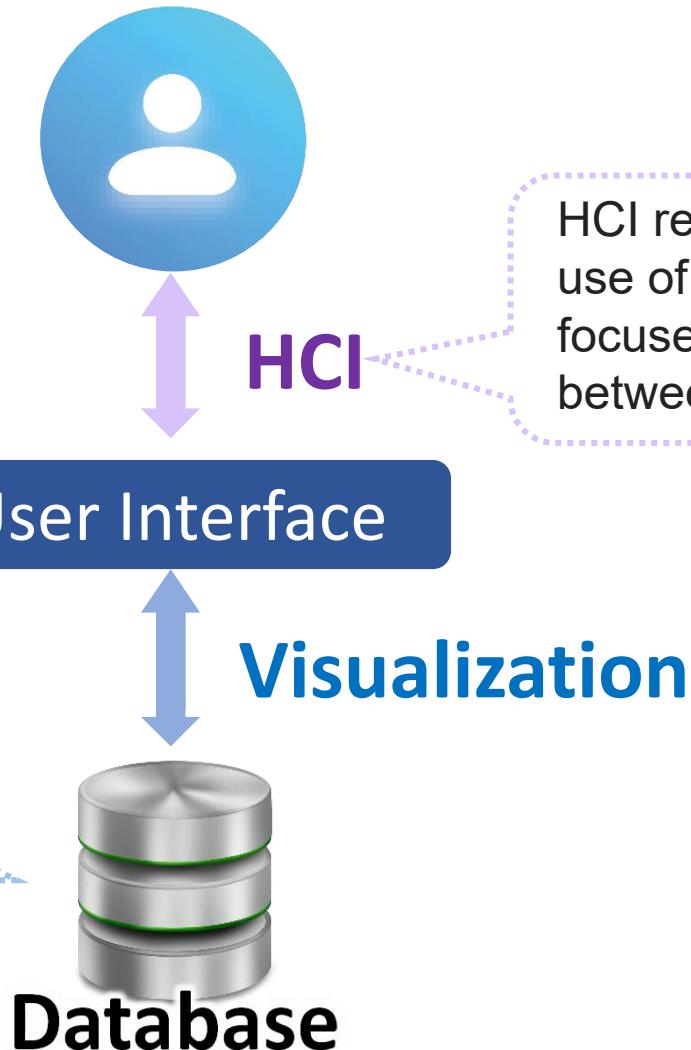
Interactive Data Exploration



Data Exploration

Data visualization refers to the techniques used to communicate data or information by encoding it as **visual objects** (e.g., points, lines or bars) contained in graphics. The goal is to communicate information **clearly and efficiently** to users.

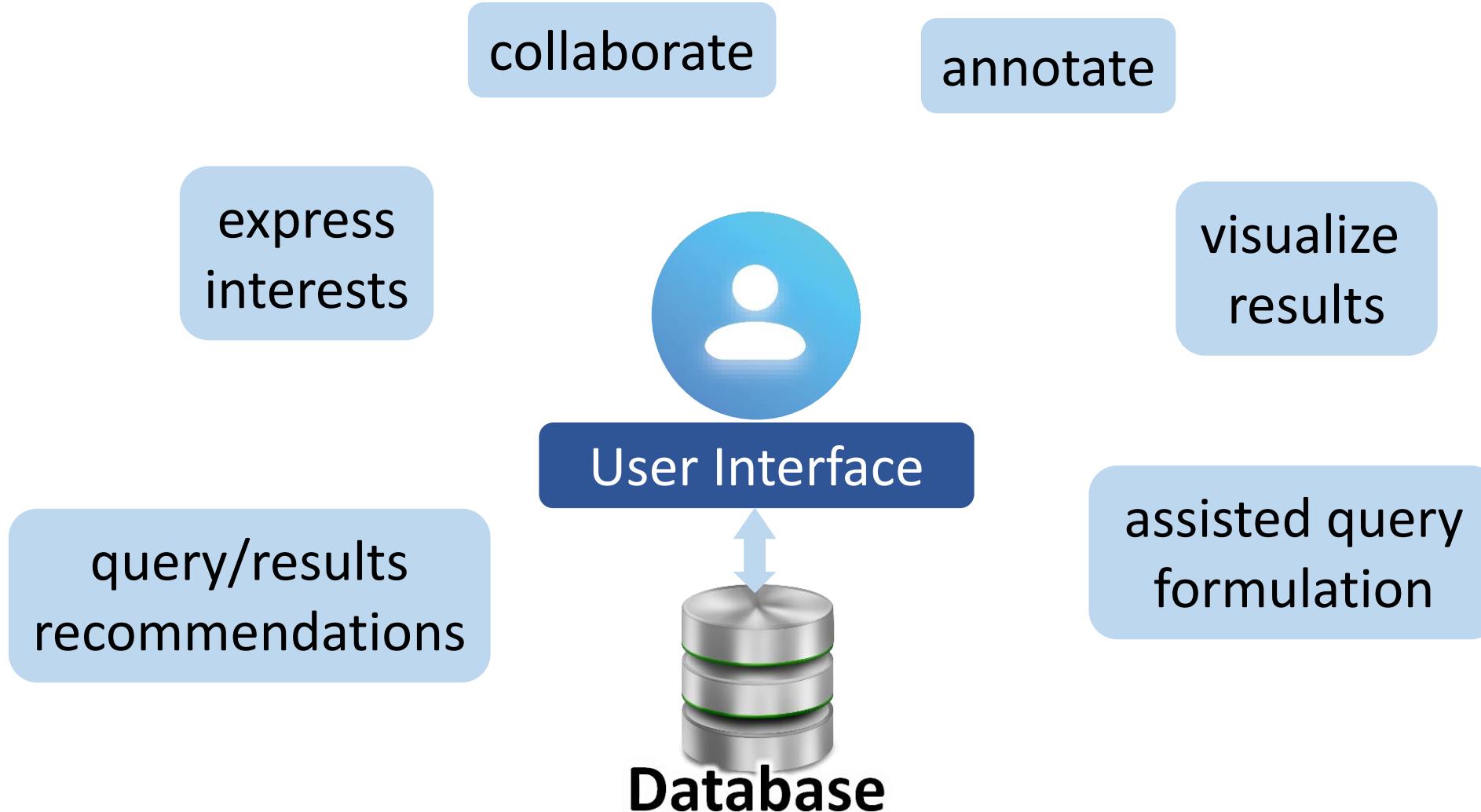
A database is an **organized collection** of data generally stored and accessed electronically from a computer system.



HCI researches the design and use of computer technology, focused on the **interfaces** between people and computers.

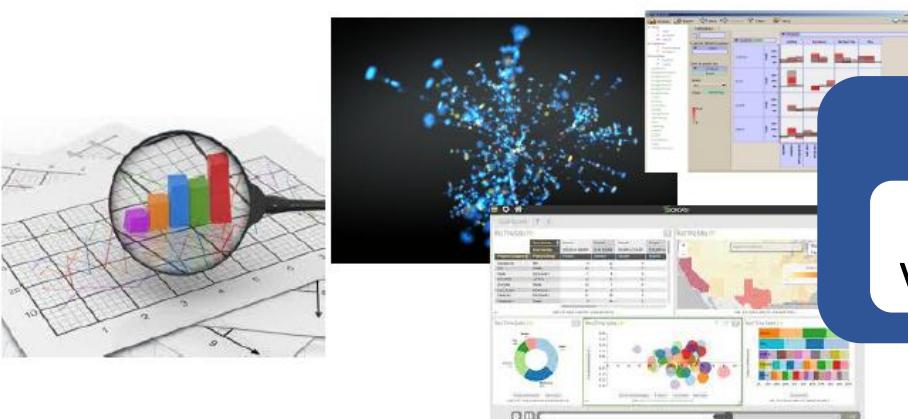


Interactive Data Exploration





Interactive Data Exploration



User Interface
Data Visualization Exploration Interface

Middleware Techniques



Database

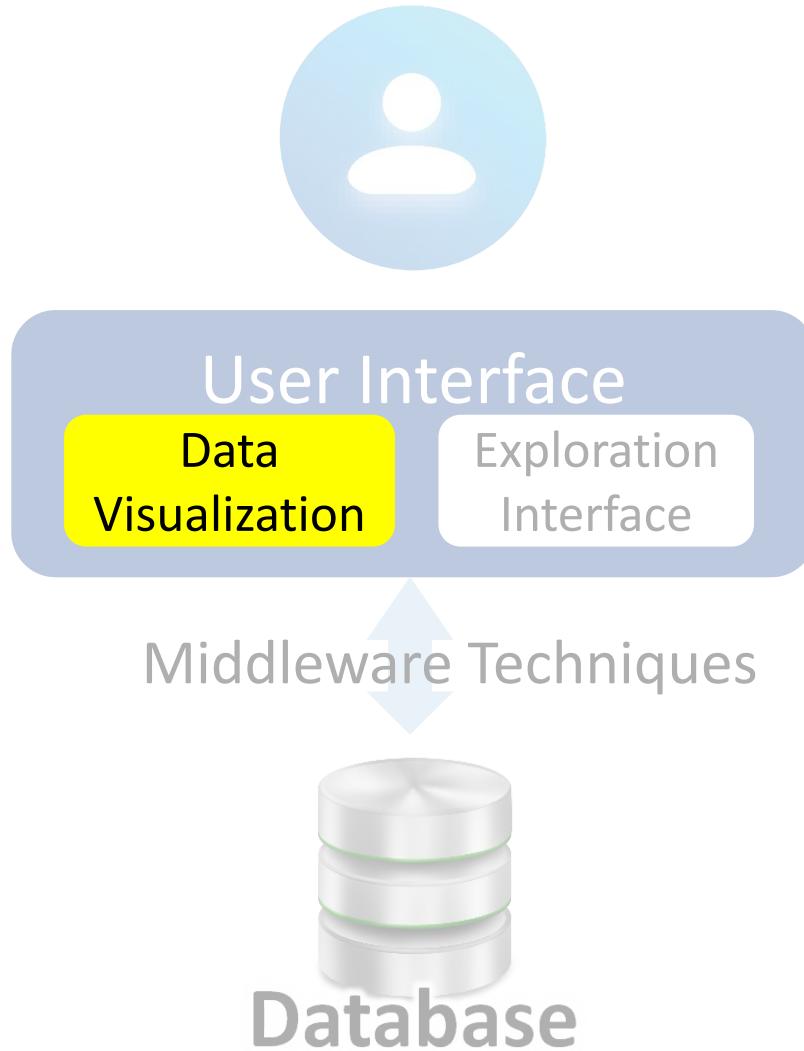




IDE: Data Visualization



Visualization Tools
Visual Optimizations
Automatic Visualization

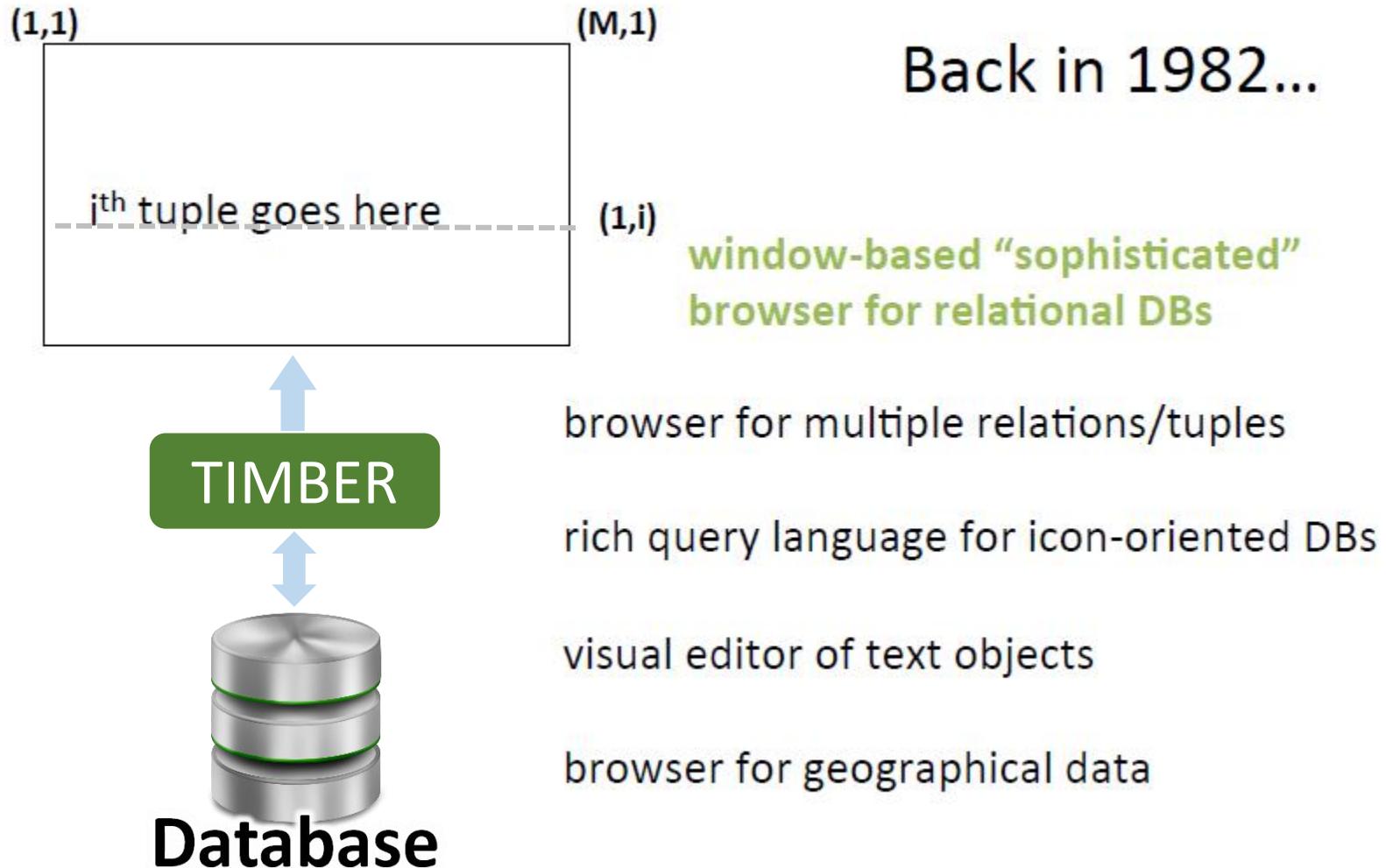


VISUALIZATION TOOLS

TIMBER

POLARIS

AstroShelf





For example, suppose an EMP relation is declared as follows:

```
CREATE EMP (name = c10, salary = i2, visual = icon)
```

Suppose the graphical token for an employee has been specified as a rectangle. It might be represented as some coding for the following information:

lines:

```
(0,0)  ->  (24,0)
(24,0) ->  (24,4)
(24,4) ->  (0,4)
(0,4)  ->  (0,0)
```

text:

```
"name ="  at (1,3)
"salary =" at (2,3)
```

fields:

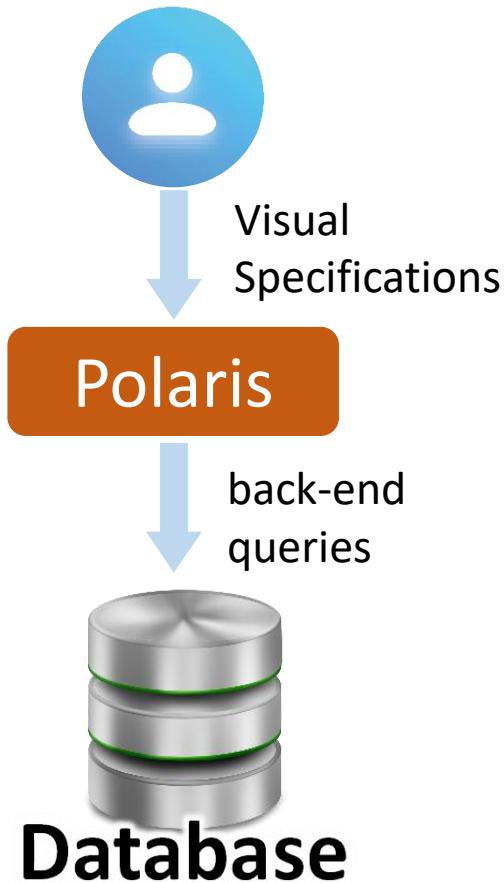
```
name    at (1,10)
salary   at (2,10)
scale-x = 1/24
scale-y = 1/24
```

and have a screen image of:

```
name = Kalash
salary = 10000
```



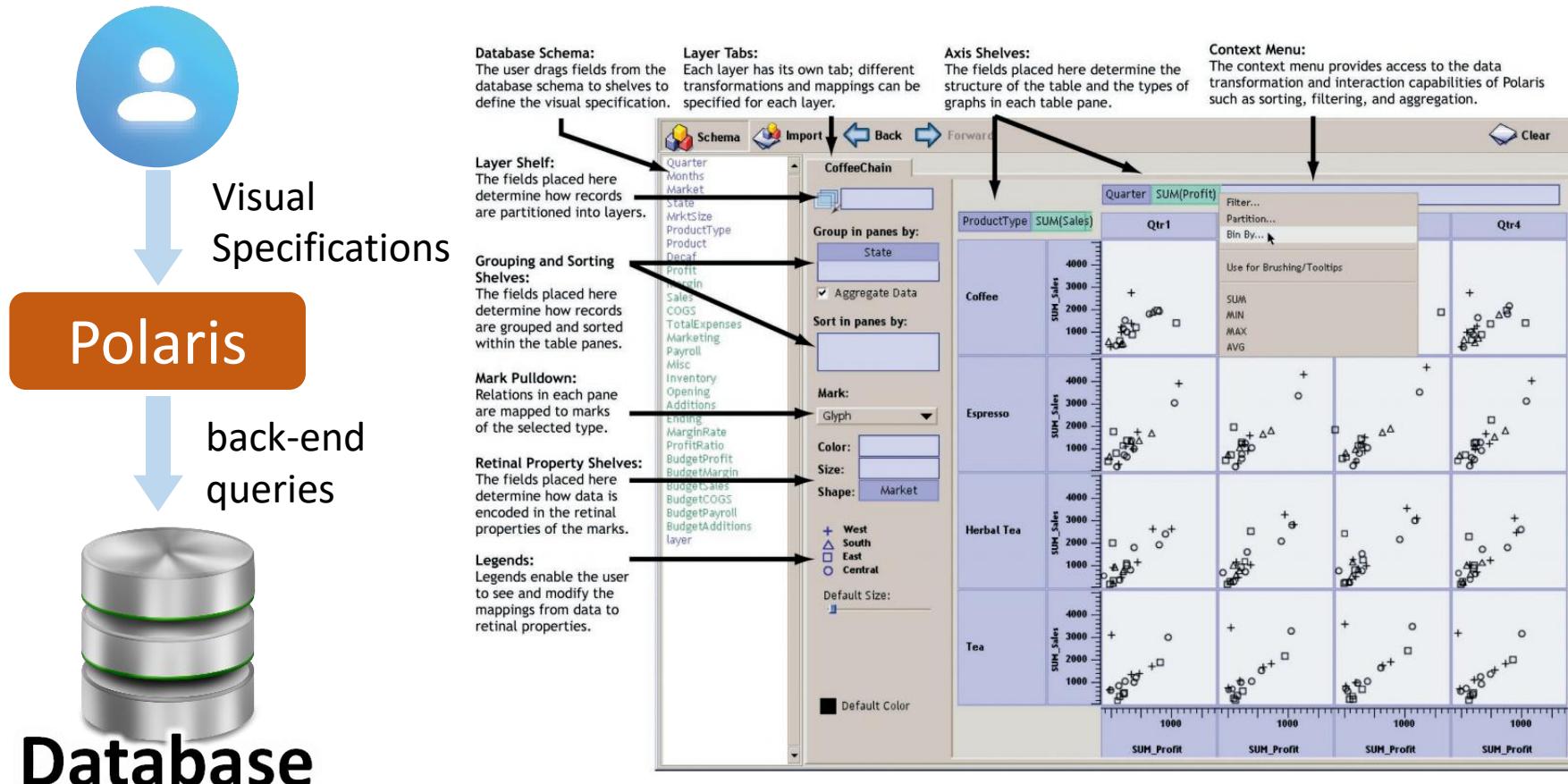
Explore large multi-dimensional databases (e.g., business, science)



- Unify tables and graphs
- Expressiveness
- Interface simplicity
- Code simplicity

	A	B	C	D	E
1	Region	East			
2					
3	Sum of Sales		Product		
4	Month	Salesperson			
5	May	Buchanan		\$17,578	\$17,578
6		Davolio	\$22,977		\$22,977
7	May Total		\$22,977	\$17,578	\$40,555
8	Jun	Buchanan	\$10,017	\$7,711	\$17,728
9		Davolio	\$6,805	\$5,575	\$12,380
10	Jun Total		\$16,822	\$13,286	\$30,108
11	Grand Total		\$39,799	\$30,864	\$70,663

Explore large multi-dimensional databases





Polaris: Operands

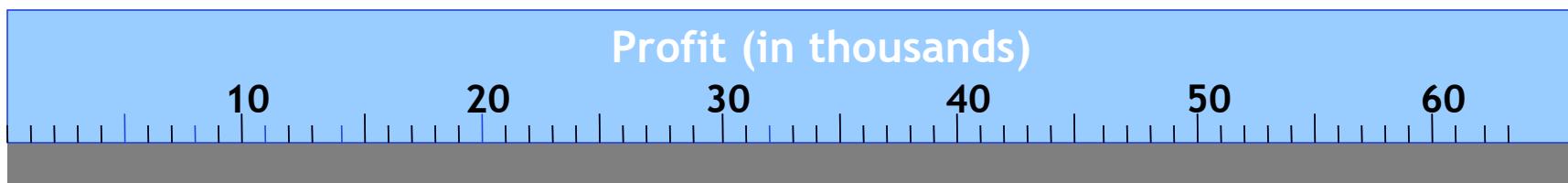
Ordinal fields - interpret domain as a set that partitions table into rows and columns:

QUARTER = {Quarter1, Quarter2, Quarter3, Quarter4} →

Quarter 1	Quarter 2	Quarter 3	Quarter 4
31,400	35,600	37,120	30,900

Quantitative fields - treat domain as single element set and encode spatially as axes:

PROFIT = {P[0 - 65,000]} →





Polaris: Concatenation Operands

Ordered union of set interpretations:

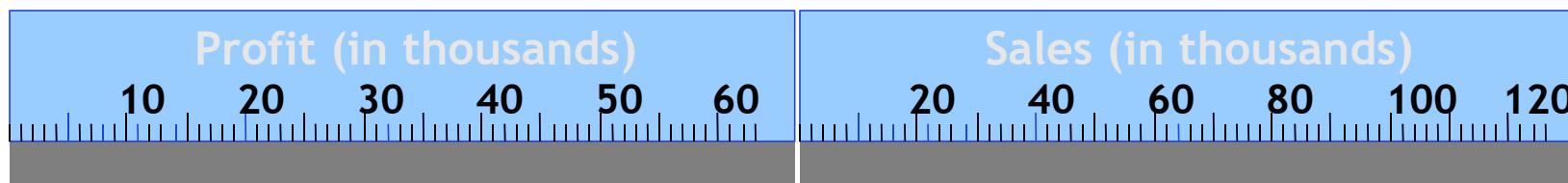
QUARTER + PRODUCT_TYPE

= {QTR1, QTR2, QTR3, QTR4} + {Coffee, Tea}

= {QTR1, QTR2, QTR3, QTR4, Coffee, Tea}

Quarter 1	Quarter 2	Quarter 3	Quarter 4	Coffee	Tea
31,400	35,600	37,120	30,900	37,120	30,900

PROFIT + SALES = {P[0-65,000], S[0-125,000]}





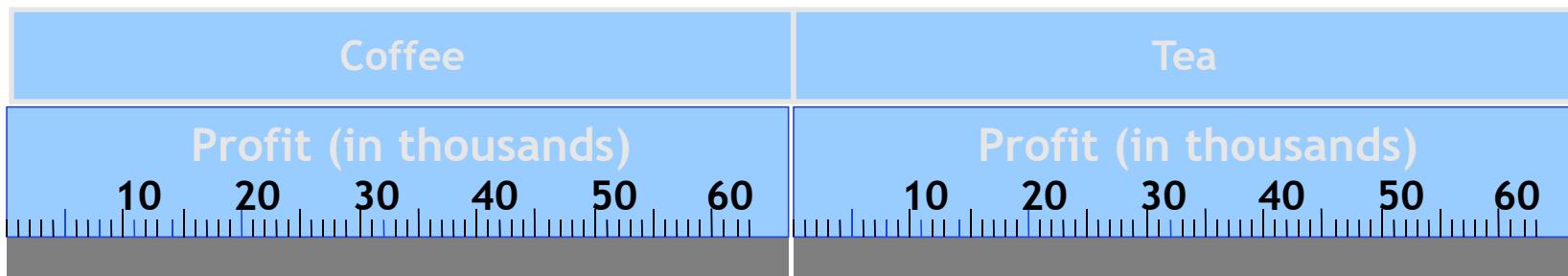
Polaris: Cross Operands

Cross-product of set interpretations:

QUARTER X PRODUCT_TYPE =

Quarter 1	Quarter 2	Quarter 3	Quarter 4	Quarter 1	Quarter 2	Quarter 3	Quarter 4
Coffee	Tea	Coffee	Tea	Coffee	Tea	Coffee	Tea
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
(Qtr1, Coffee)	(Qtr1, Tea)	(Qtr2, Coffee)	(Qtr2, Tea)	(Qtr3, Coffee)	(Qtr3, Tea)	(Qtr4, Coffee)	(Qtr4, Tea)

PRODUCT_TYPE X PROFIT =



Polaris: Nest Operands



QUARTER X MONTH

would create entry twelve entries for each quarter i.e. (Qtr1, December)

QUARTER / MONTH

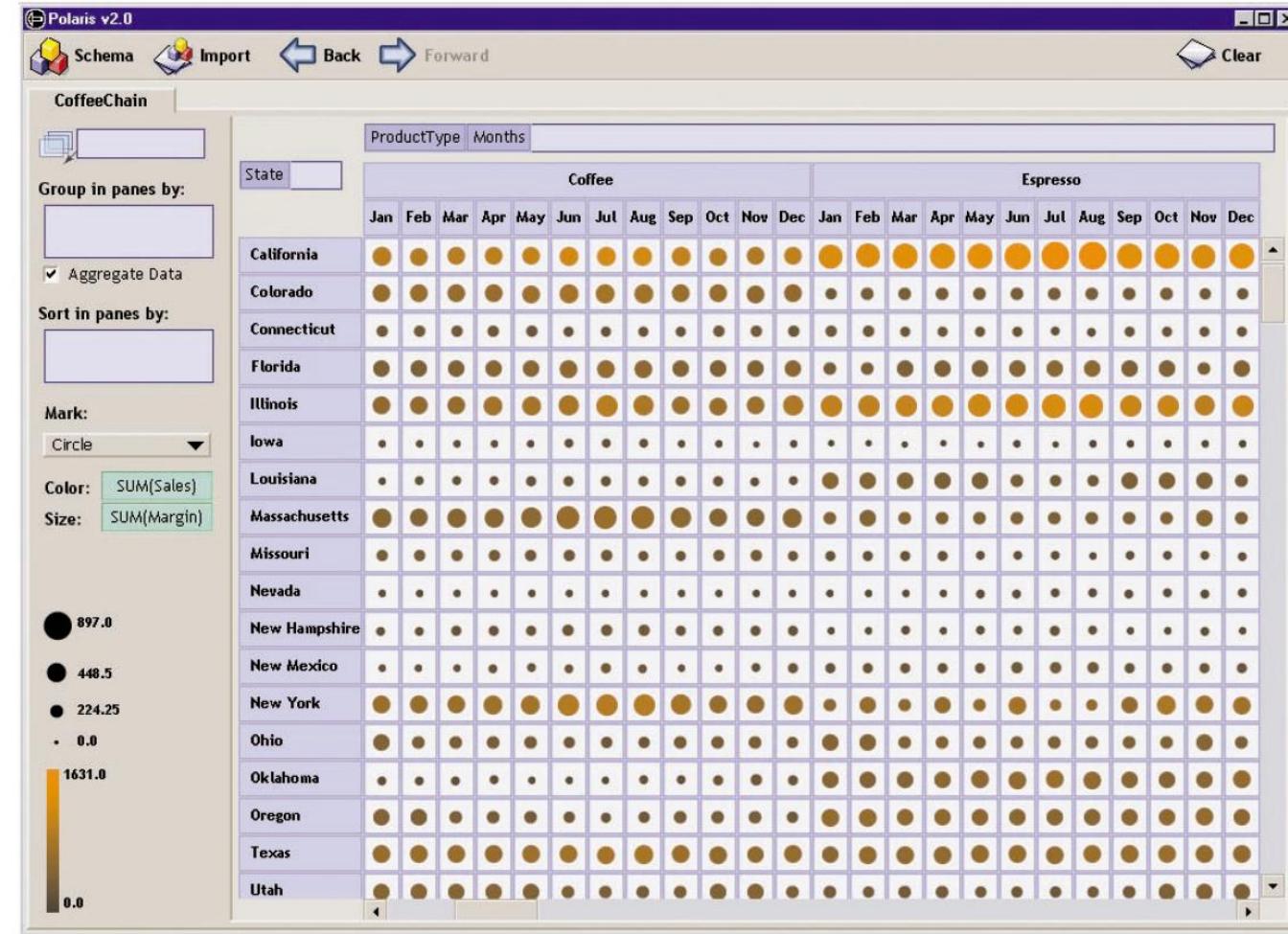
would only create three entries per quarter

Qtr1			Qtr2			Qtr3			Qtr4		
Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec



Polaris: formalism

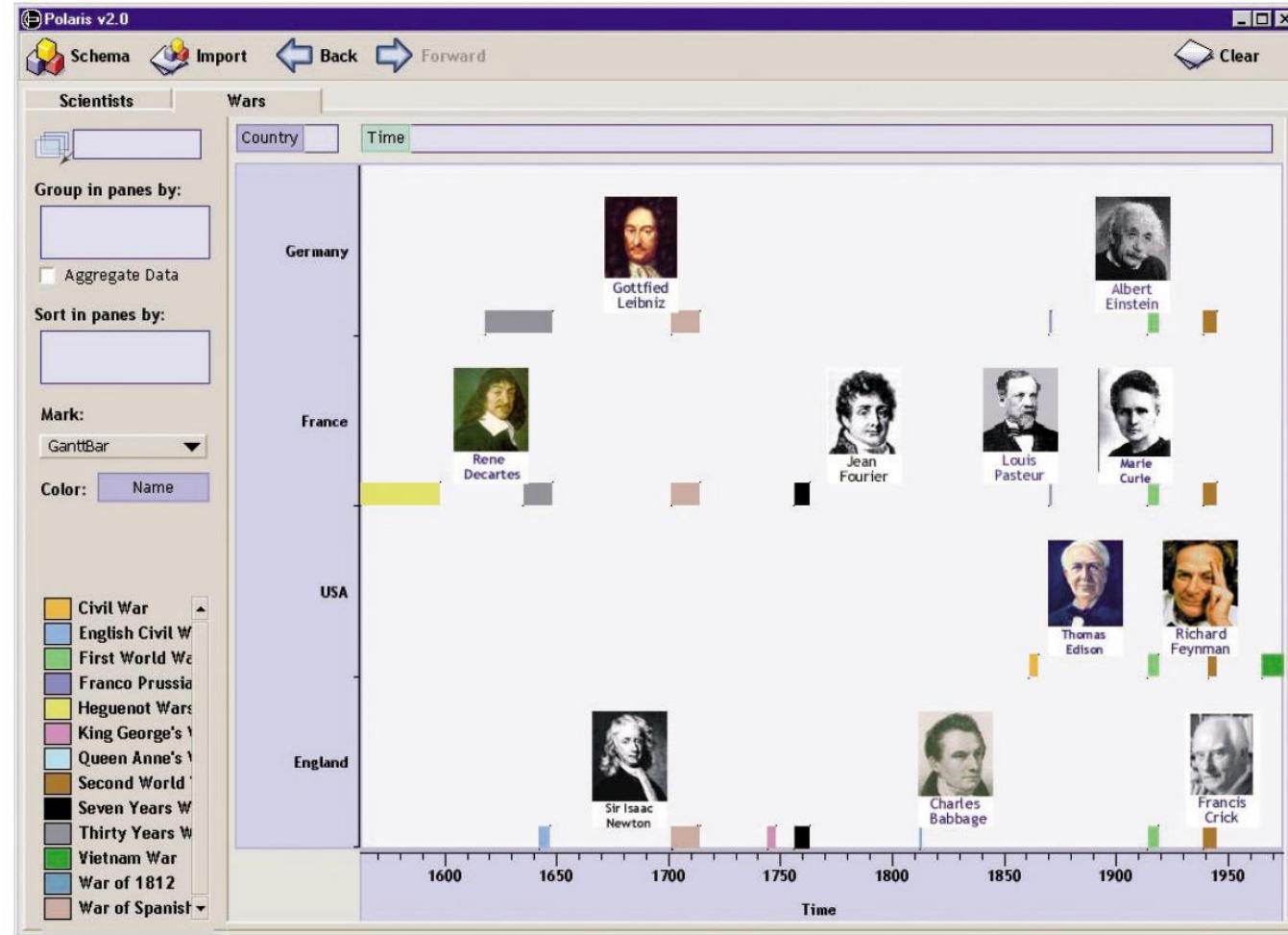
Specification of different graph types





Polaris: formalism

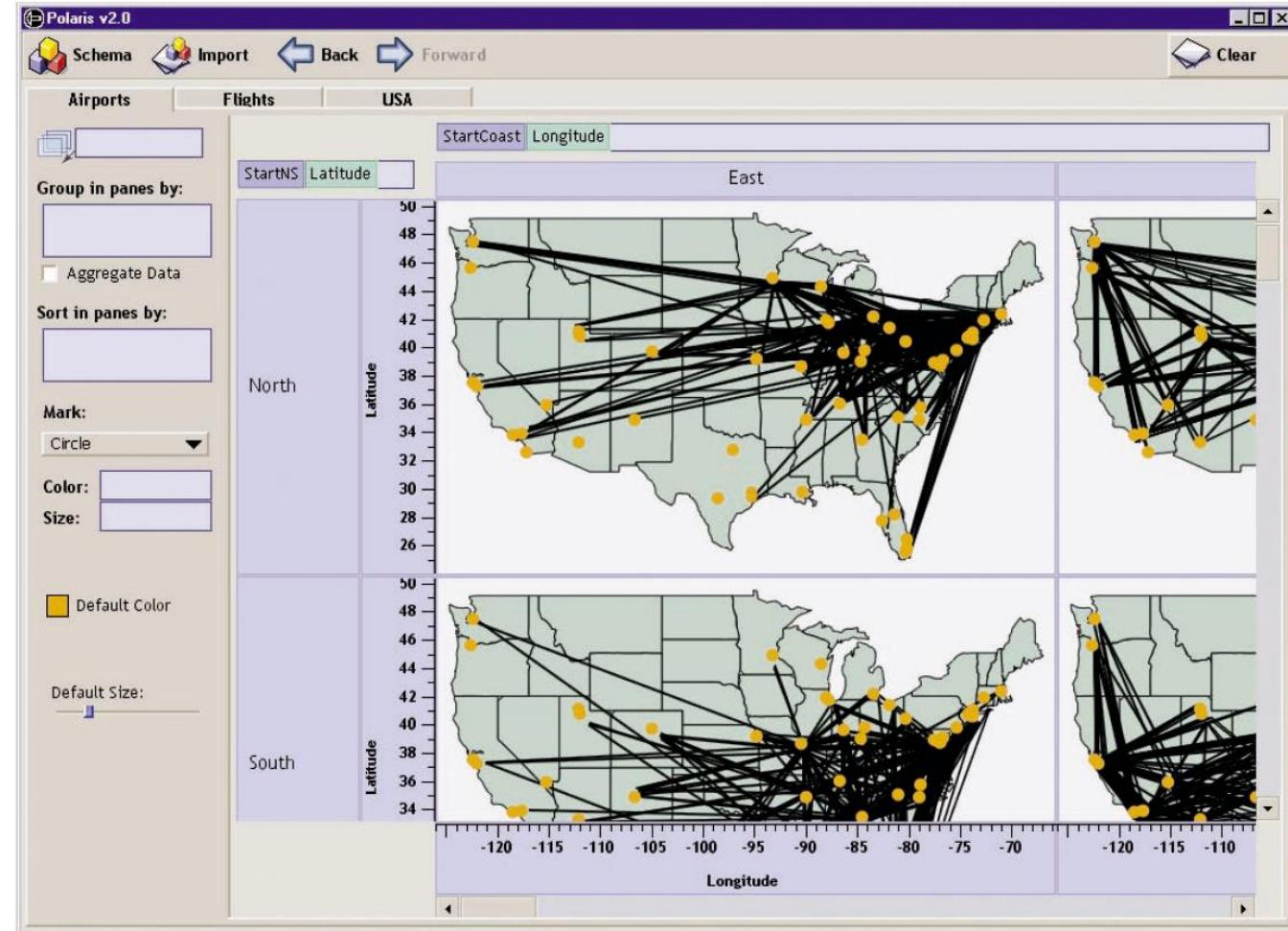
Specification of different graph types



Polaris: formalism



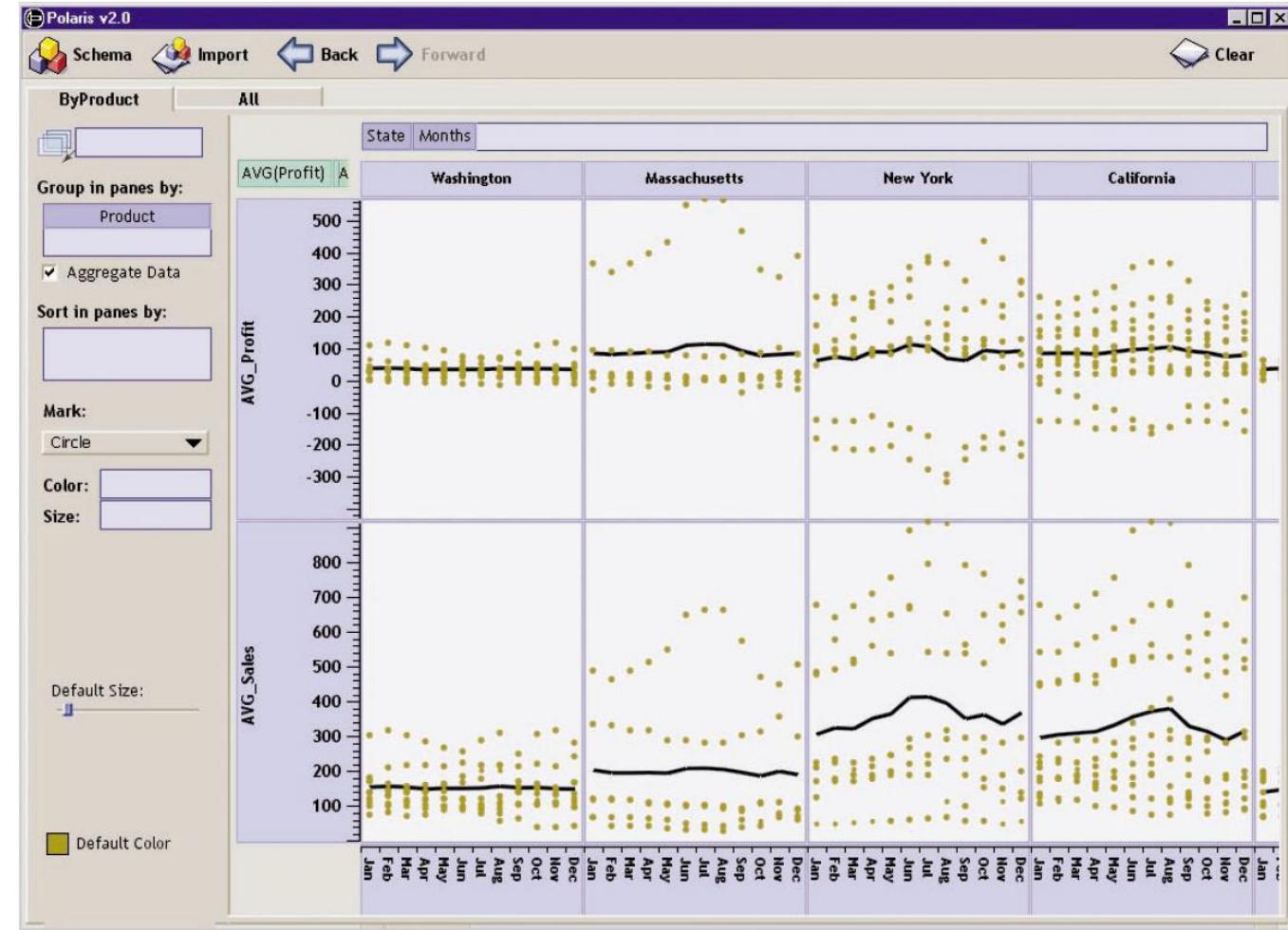
Specification of different graph types





Polaris: formalism

Specification of different graph types





Polaris: formalism

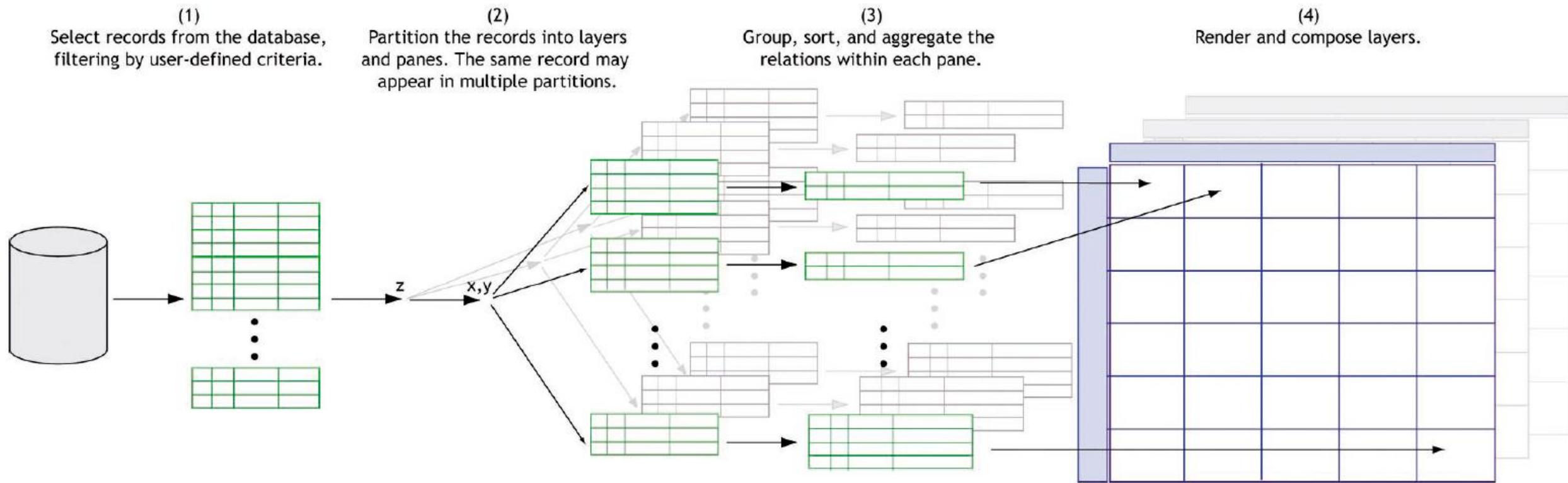
Encoding of data as retinal properties of marks in graphs

property	marks	ordinal/nominal mapping	quantitative mapping
shape	glyph	○ □ + △ S U	
size	rectangle, circle, glyph, text	• ● ○ ○ ○	• • • • • • • • •
orientation	rectangle, line, text	— — / \ \	— - - - / / / / / /
color	rectangle, circle, line, glyph, y-bar, x-bar, text, gantt bar	■ ■ ■ ■ ■ ■ ■ ■ ■ ■ ...	min max A horizontal color bar showing a gradient from dark gray on the left to red on the right, with the words "min" and "max" at the ends.



Polaris: formalism

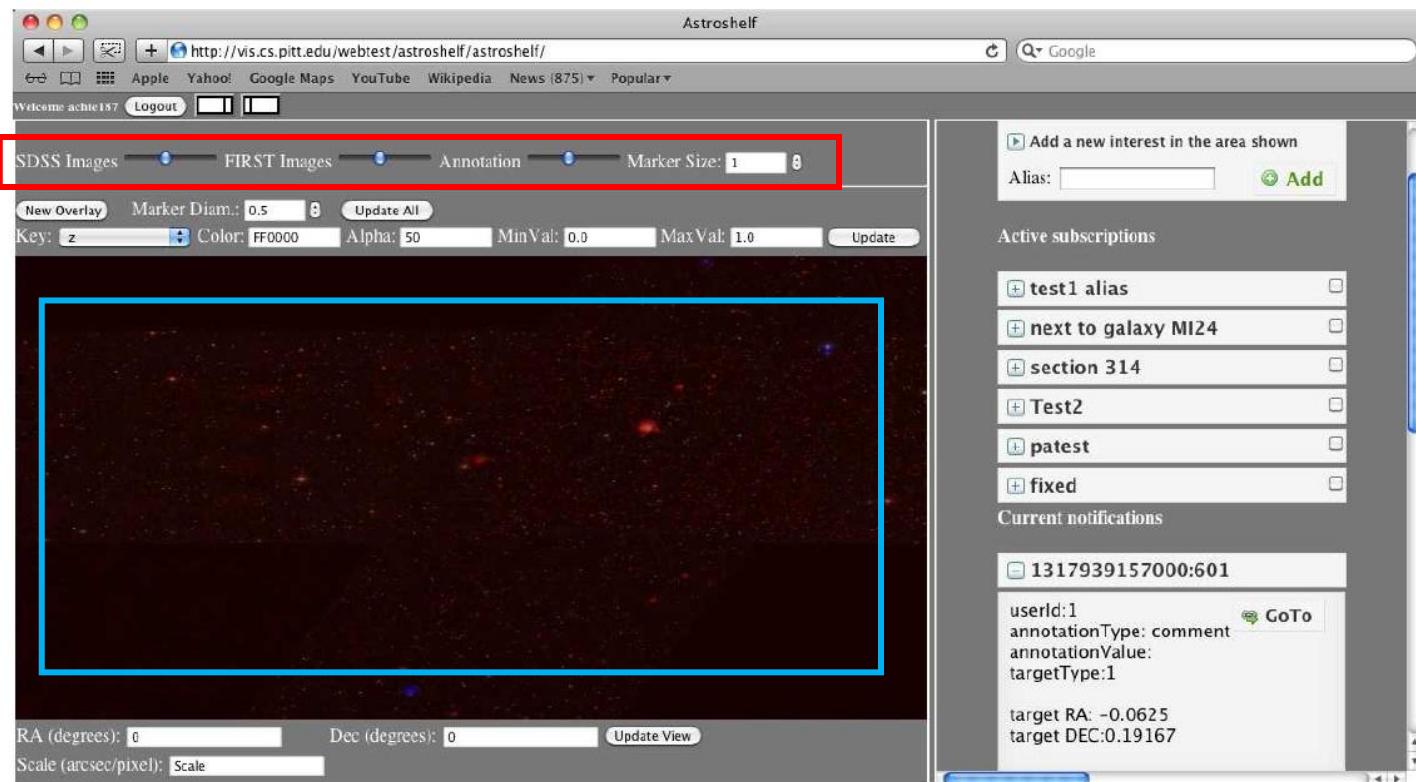
Translation of visual specification into SQL queries





AstroShelf: goals

- Scalable annotation framework
- Continuous workflow enactment system





AstroShelf



The screenshot shows the AstroShelf web application interface. On the left, a search panel titled "Search Objects" contains the following fields:

- Guided SQL**, **Direct SQL**, **Query History** buttons.
- Surveys:** SDSS FIRST
- From:** tables: PhotoObj as p LEFT OU ▾
- Parameters:** objid RA,Dec type
name u, err_u g, err_g
r, err_r i, err_i z, err_z
redshift, red_err specclass zconf
ALL NONE DEFAULT
Add Parameters button.
- Conditions:** (empty)
- Limitation:** (empty)
- SQL Query:** enable user preference
- Buttons:** Search, Cancel, Reset.

At the bottom of the search panel, it displays: RA,Dec (degrees): 1.70000000, 1.05500000 – Scale (arcsec/pixel): 0.10.

The main area features a dark background with numerous small white stars of varying sizes, representing a celestial map. To the right of the map, there is a vertical sidebar with several tabs: Overlays, Results, Object details, Trend Image, Thumbnails, and a currently selected tab labeled Thumbs.



AstroShelf



University of Pittsburgh
NSF award OIA-1028162

Computer Science:
Alexandros Labrinidis
Panos Chrysanthis
Liz Marai

Astronomy:
Jeff Newman
Michael Wood-Vasey
Arthur Kosowsky

Narrated by Eric Gratta

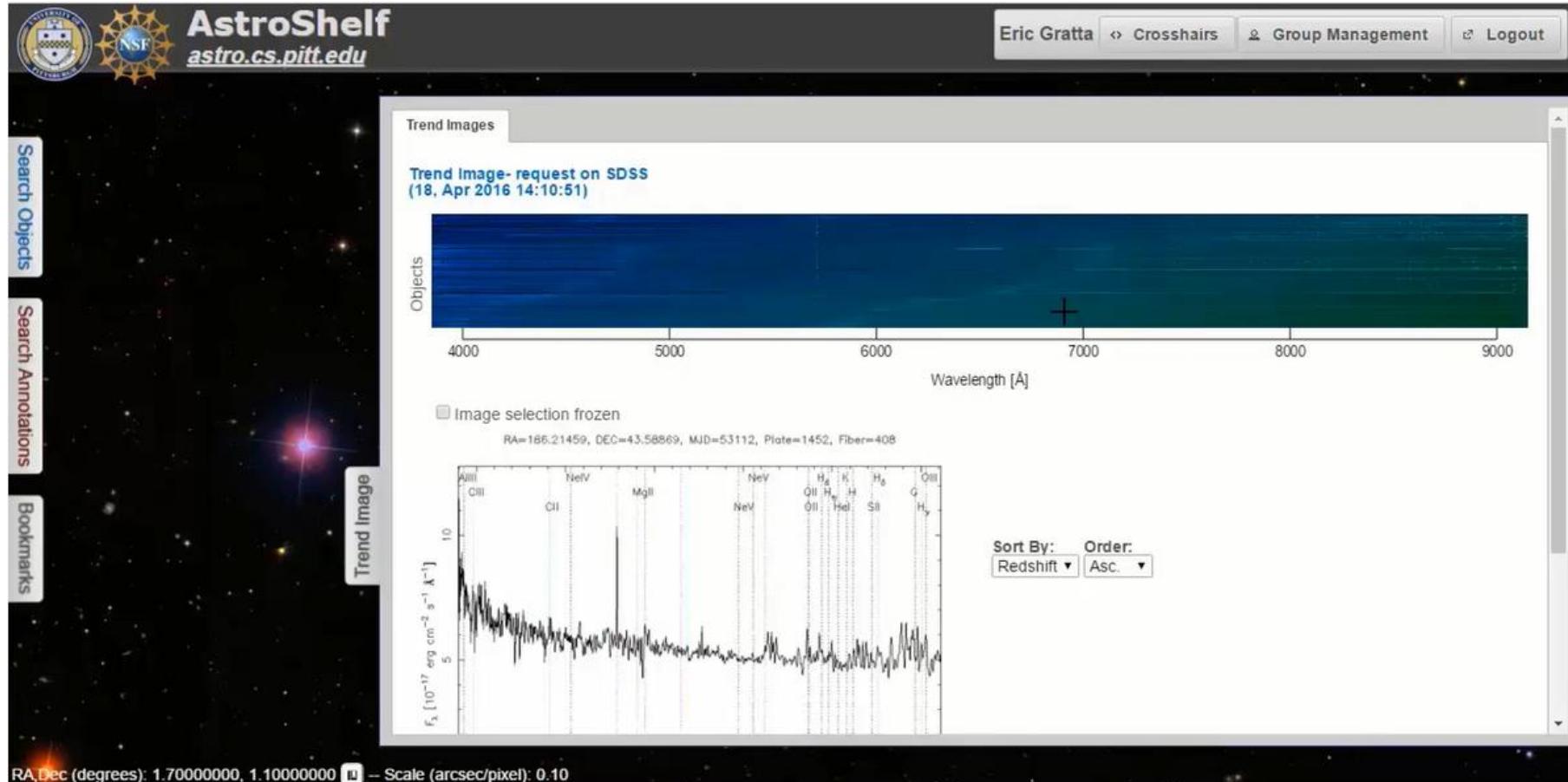
The image shows the AstroShelf interface. On the left, there's a sidebar with buttons for "Search Objects", "Search Annotations", and "Bookmarks". The main area displays a star field with several highlighted objects. A prominent orange button at the bottom left says "Results". At the top right, there's a user profile for "Eric Gratta" with options for "Crosshairs", "Group Management", and "Logout". Below the user profile is a table titled "Search Objects - request on SDSS (18, Apr 2016 13:51:13)". The table has columns for "objid", "ra", "dec", and "Object details". It lists 10 entries from 100 total, with each entry having a "more" link. At the bottom of the table, there are links for "Hide/Show", "Remove from the list", "Create overlay", "Thumbnails", "Trend Image", and "Store Result".

objid	ra	dec	Object details
587728670413750443	117.32633182	37.33628748	more
587722984425259161	181.85645294	1.03191433	more
587741421100925212	130.38813832	20.09050101	more
587726031726837811	214.64377786	1.5032223	more
588017949889331343	186.21458746	43.58868918	more
588011122504499423	192.06441827	61.46900267	more
587724241228857491	46.88097501	-7.58226202	more
587725503949832426	260.99210796	52.92284143	more
587735744229736880	247.11170566	27.98910951	more
588010879305056615	200.90333626	4.82510721	more





AstroShelf



Trend images are used to visualize overall patterns in a data set.

A screenshot of the AstroShelf application interface. At the top, there is a header bar with the university's logo, the NSF logo, the text 'AstroShelf', and the URL 'astro.cs.pitt.edu'. On the right side of the header, there are user profile links for 'Eric Gratta' and options for 'Crosshairs', 'Group Management', and 'Logout'. The main area shows a dark background with numerous small stars and a few larger, more prominent ones. A central modal dialog box is open, titled 'Annotate'. It contains fields for 'Title' (with a dropdown menu showing 'Shared', 'Private', 'Shared' (which is selected), and 'Public'), a section for selecting an annotation type ('Text', 'Tag', or 'Link'), and a 'Tags/Keywords' input field. At the bottom of the dialog are 'Ok' and 'Cancel' buttons. Along the left and right edges of the main window, there are vertical menus with items like 'Search Objects', 'Search Annotations', 'Bookmarks', 'Overlays', 'Results', 'Object details', 'Trend Image', and 'Thumbnails'. At the bottom of the screen, there is a status bar with the text 'RA,Dec (degrees): 186.21458746, 43.60368918' and a scale indicator 'Scale (arcsec/pixel): 0.10'.

Annotate

Title:

Select an annotation type and fill out the form:

Text Tag Link

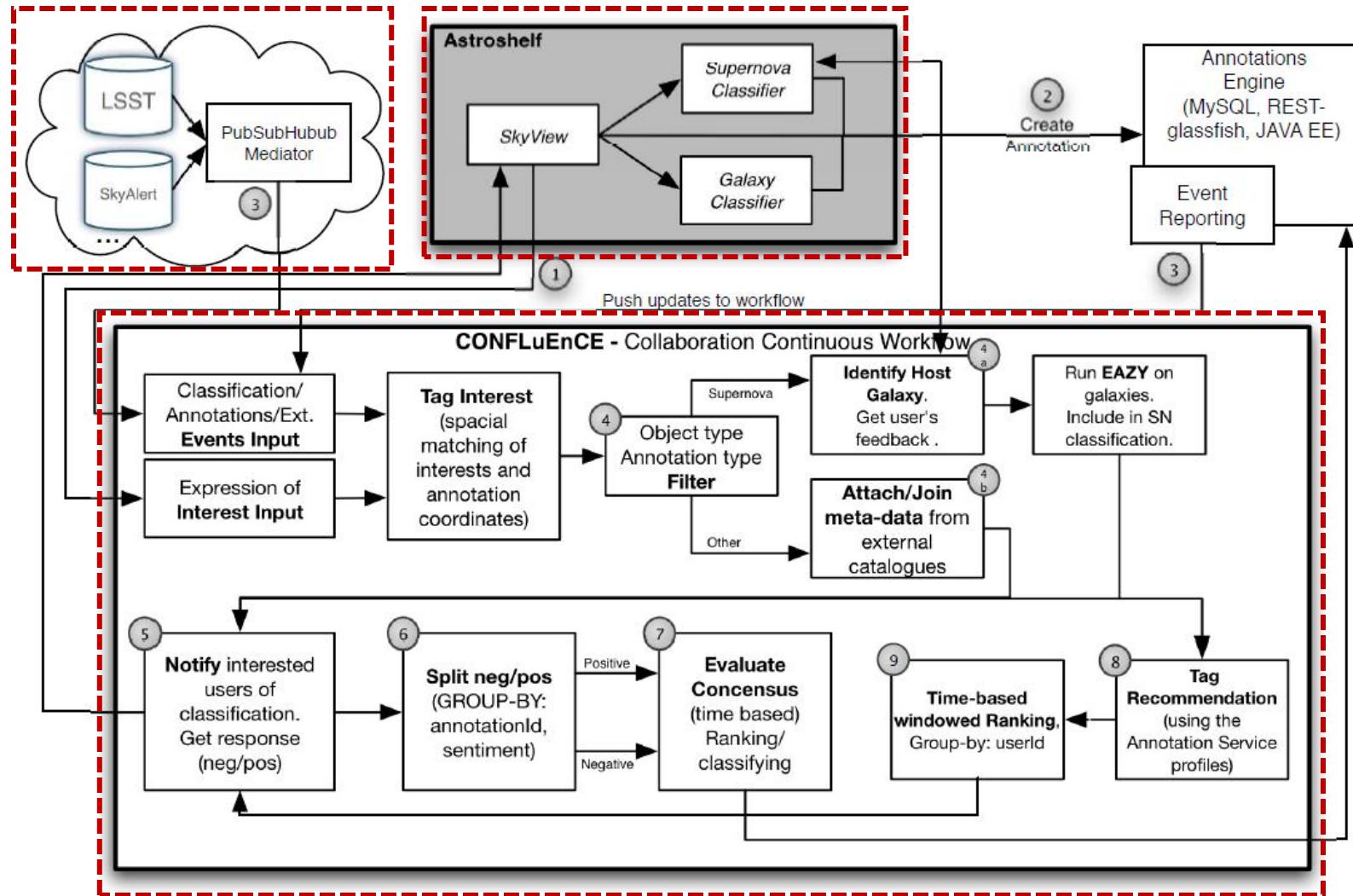
Tags/Keywords

Ok Cancel

RA,Dec (degrees): 186.21458746, 43.60368918 – Scale (arcsec/pixel): 0.10

Annotations: <annotation type, annotation value>.

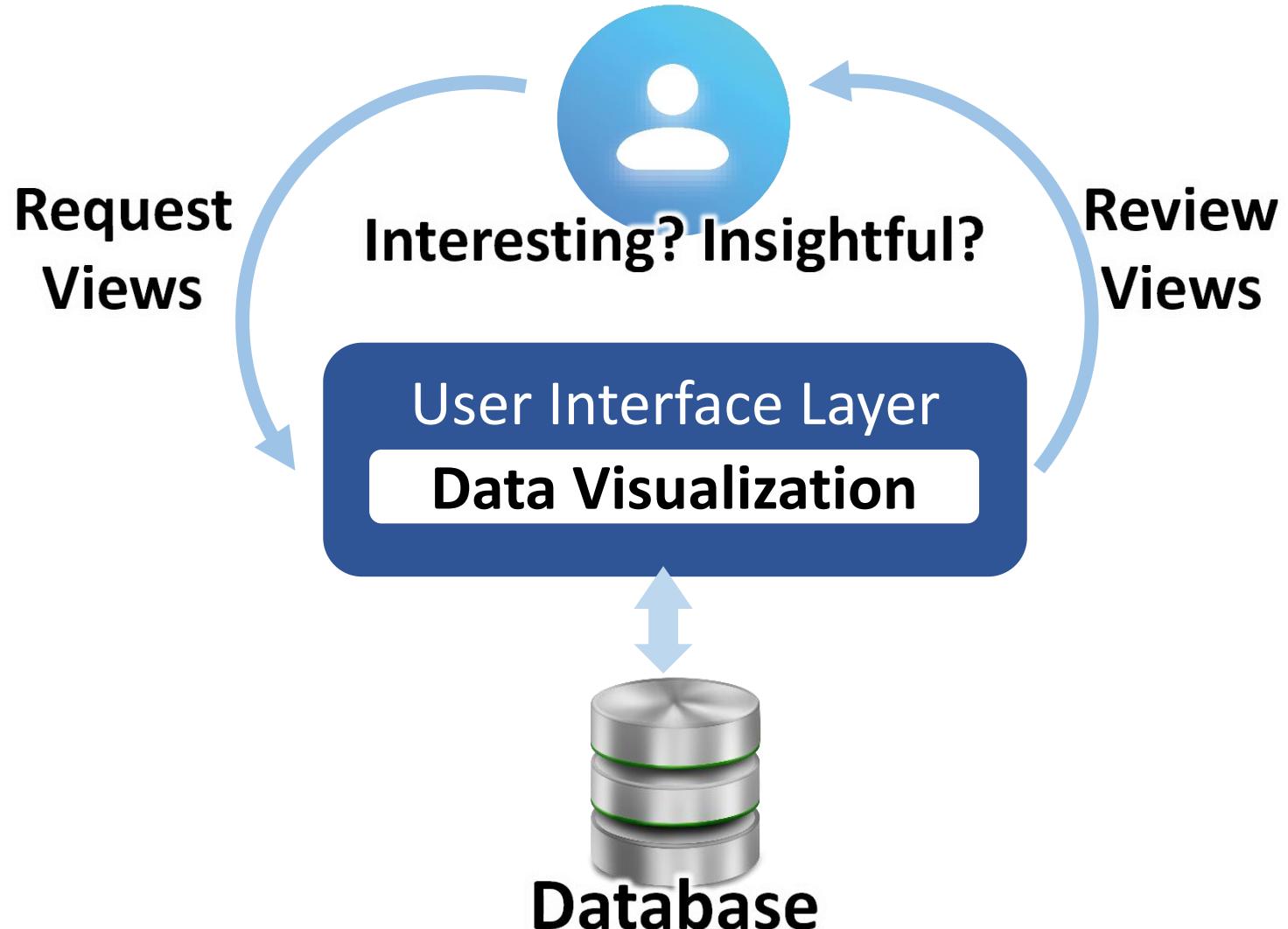
AstroShelf

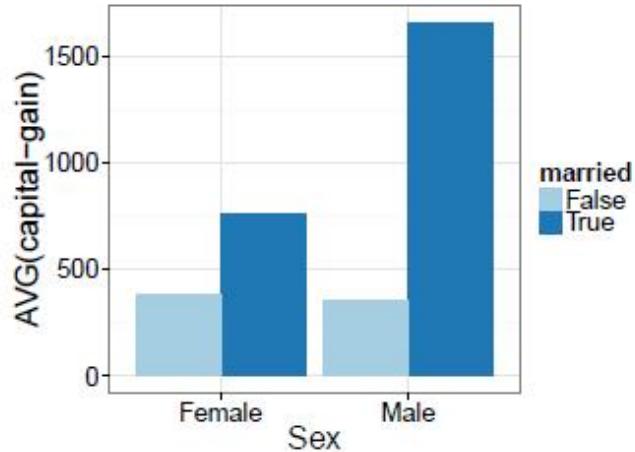




Automatic Visualization

manual, repetitive exploration for best visualization(s)

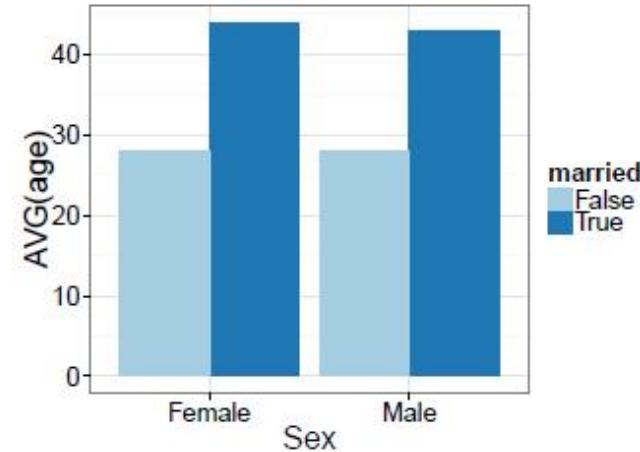




(a) Interesting Visualization

Sex	Avg Capital Gain
Unmarried Adults	
Female	380
Male	356
Married Adults	
Female	758
Male	1657

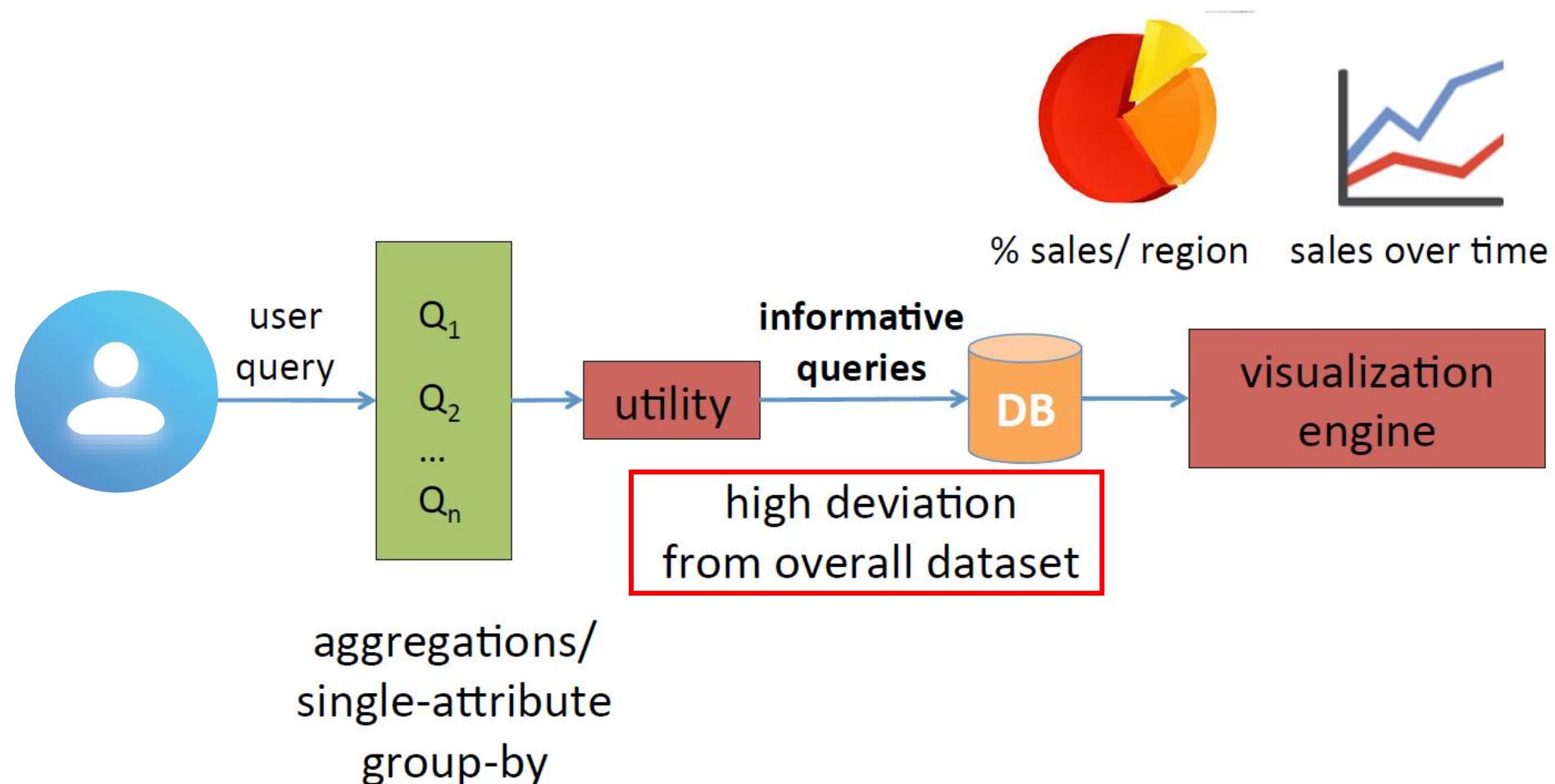
(c) Data: Avg Capital Gain vs. Sex



(b) Uninteresting Visualization

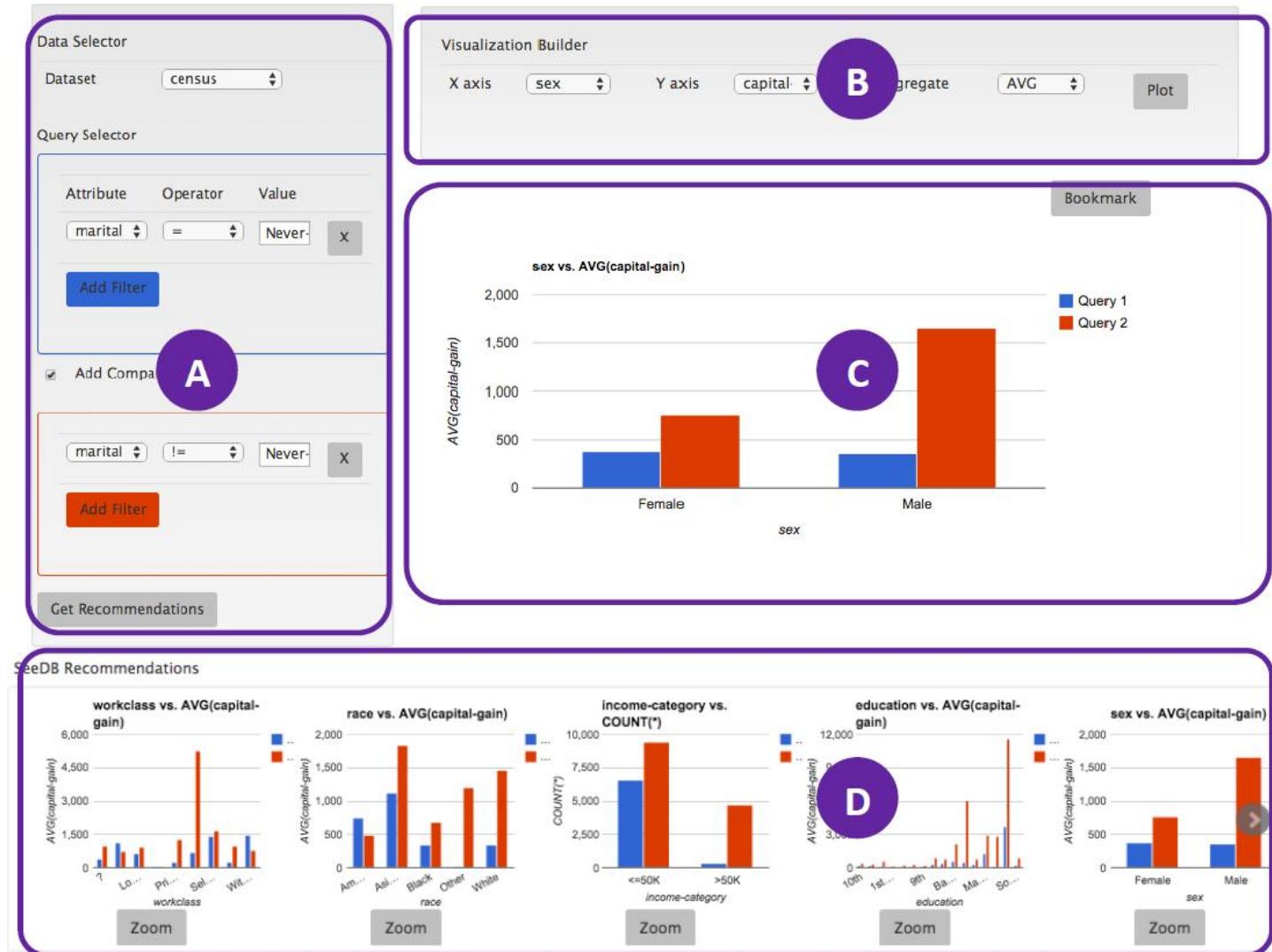
Sex	Avg Age
Unmarried Adults	
Female	28
Male	28
Married Adults	
Female	44
Male	43

(d) Data: Avg Age vs. Sex





SeeDB



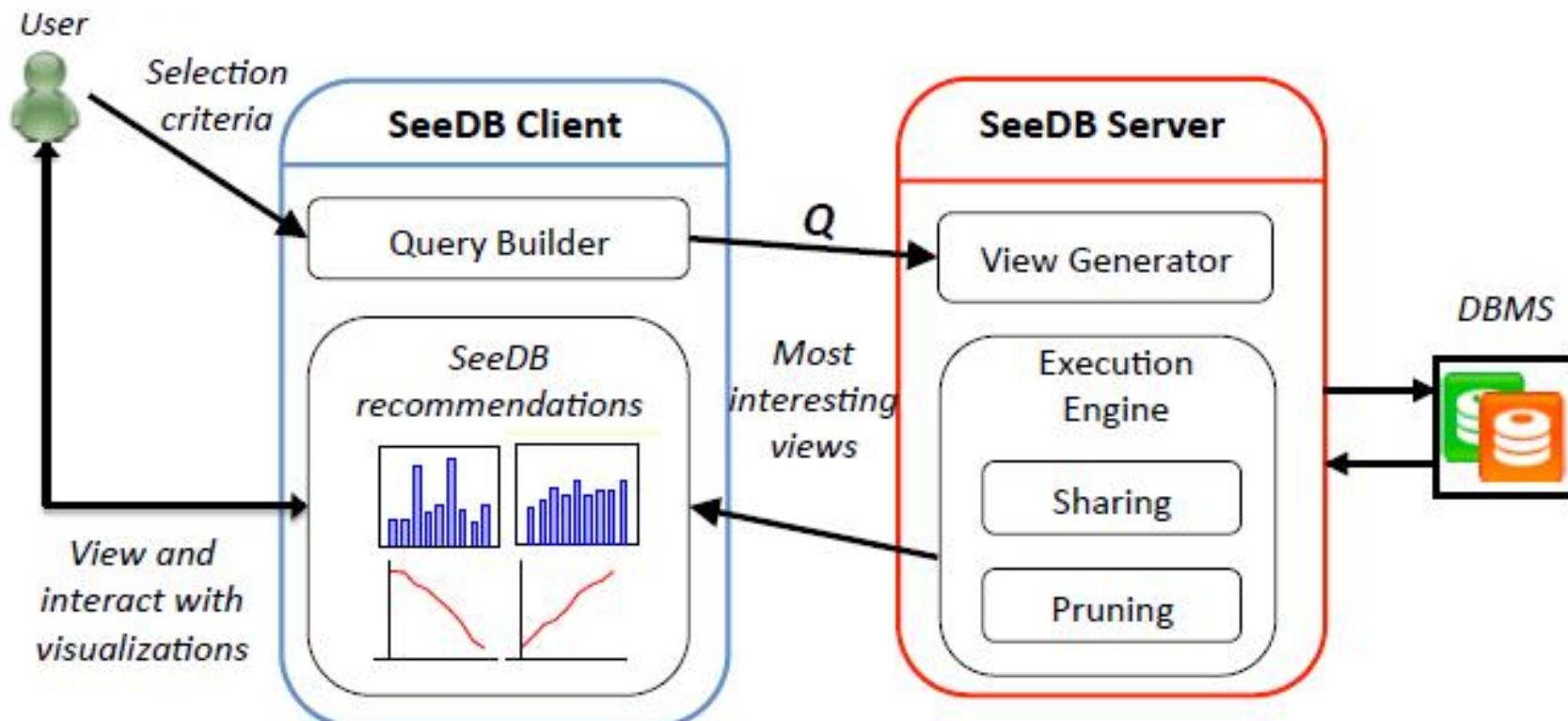


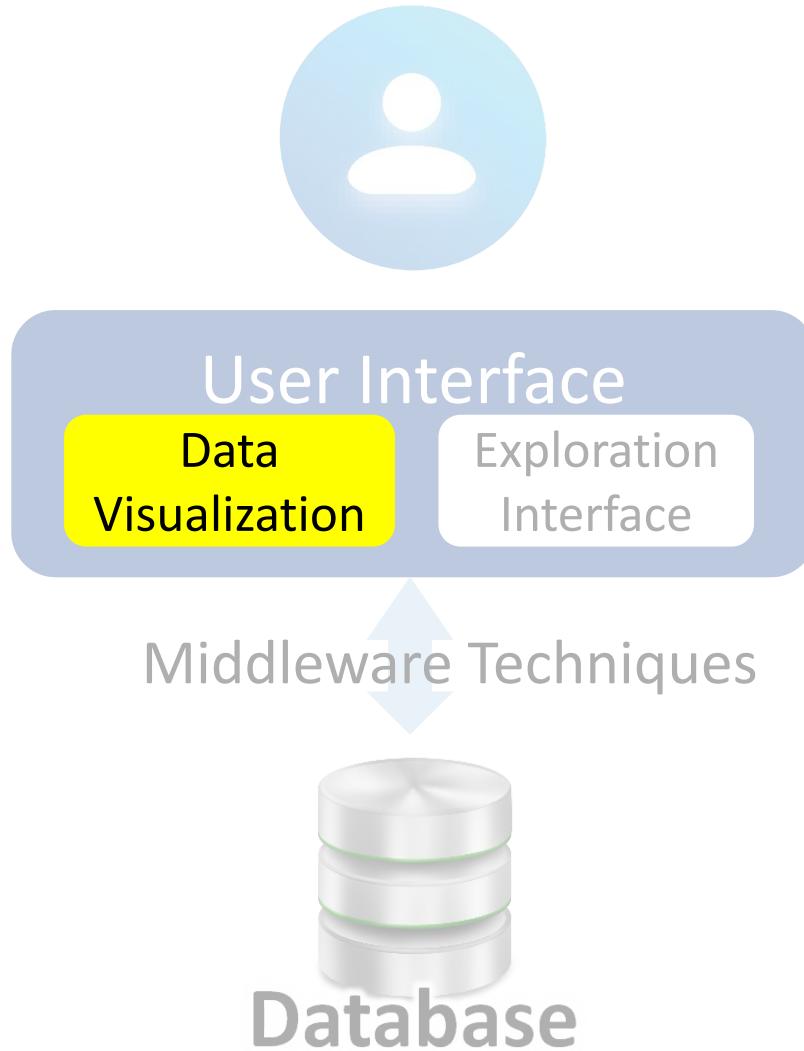
Figure 3: SeeDB Architecture



IDE: Data Visualization



Visualization Tools
Visual Optimizations
Automatic Visualization



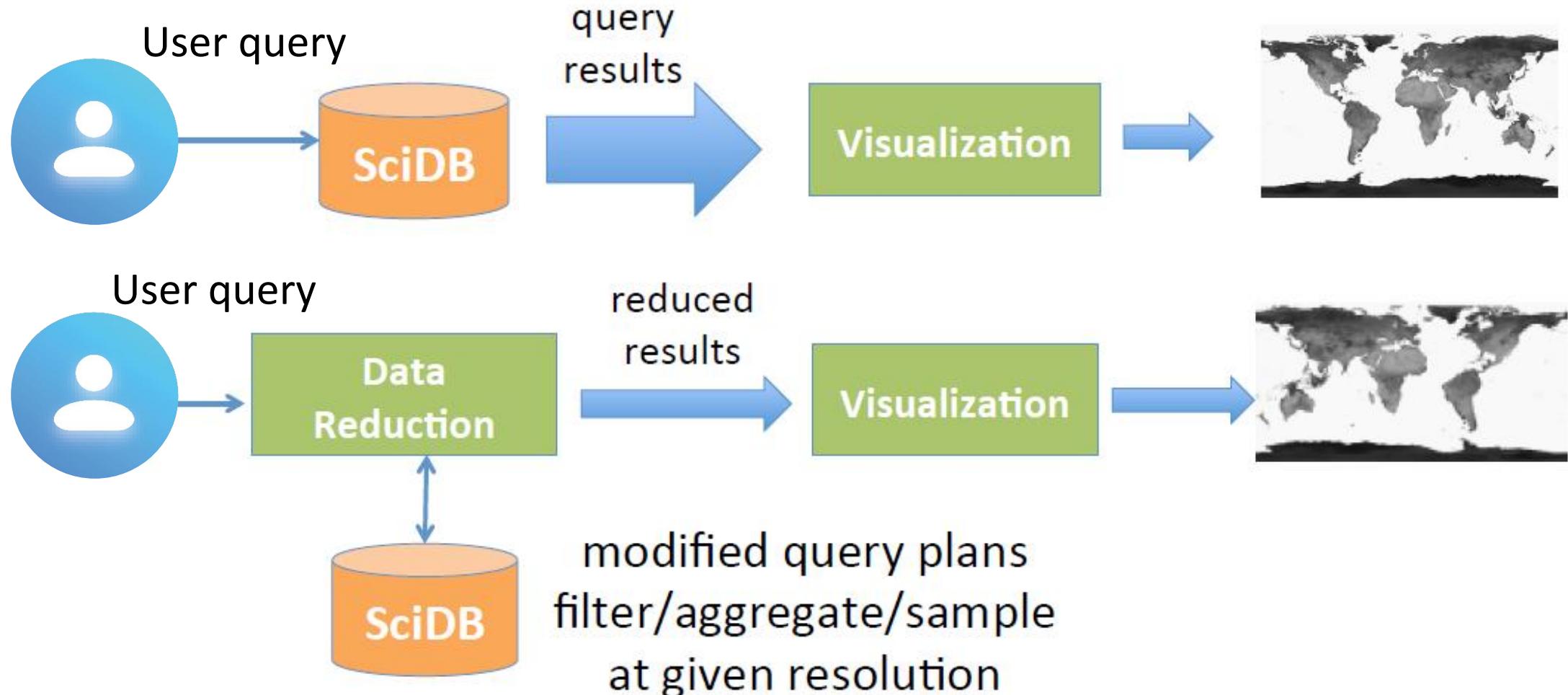
探索式数据分析跟传统数据分析相比有什么不同？

- A 探索式数据分析不涉及数据库操作
- B 传统数据分析不涉及可视化
- C 探索式数据分析中用户难以通过直接查询获得想要分析的信息

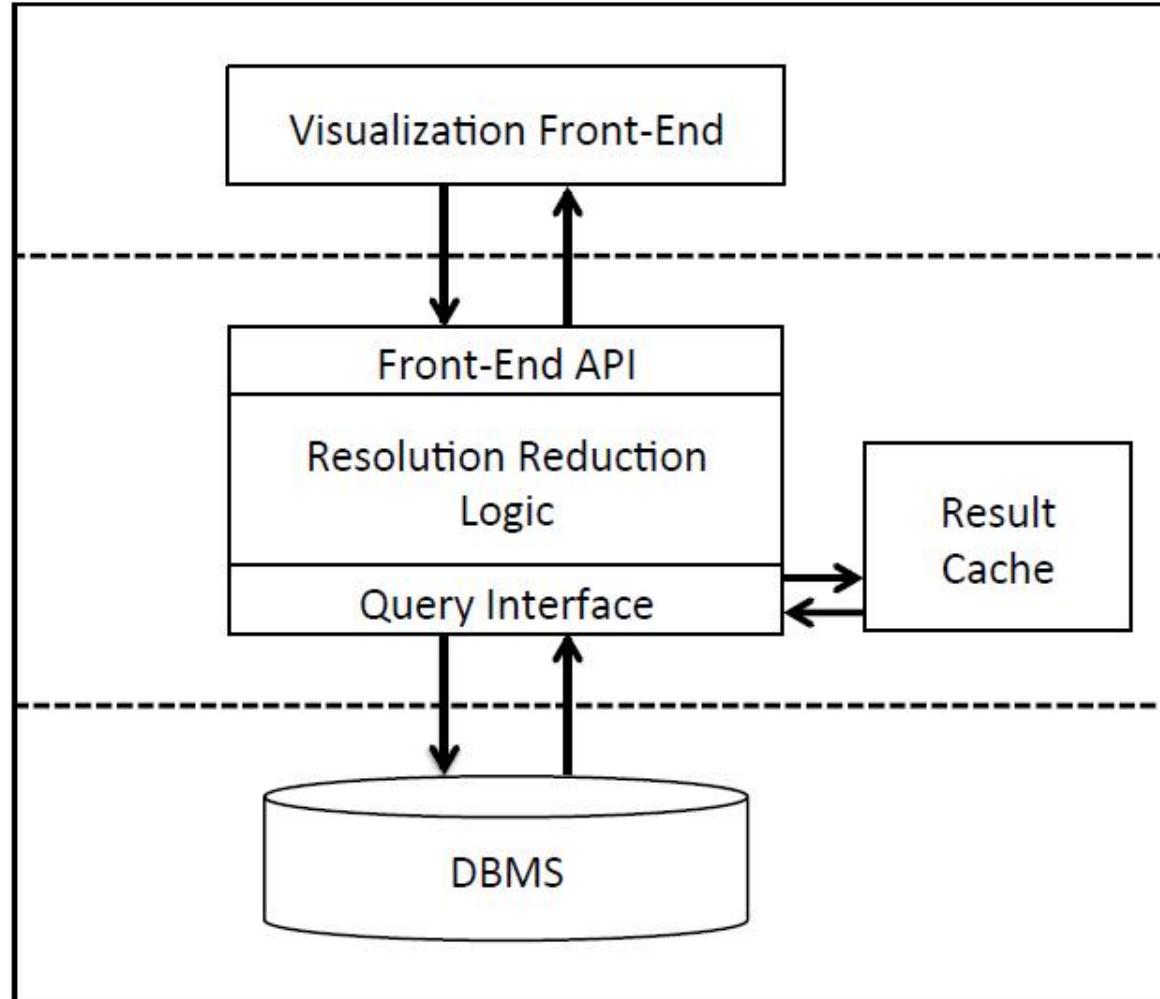
提交

Visual Optimization

ScalaR: Resolution Reduction



ScalaR: resolution reduction

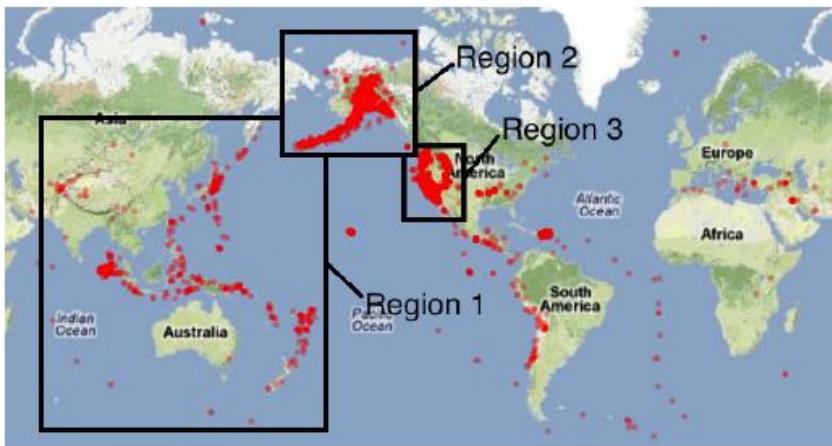




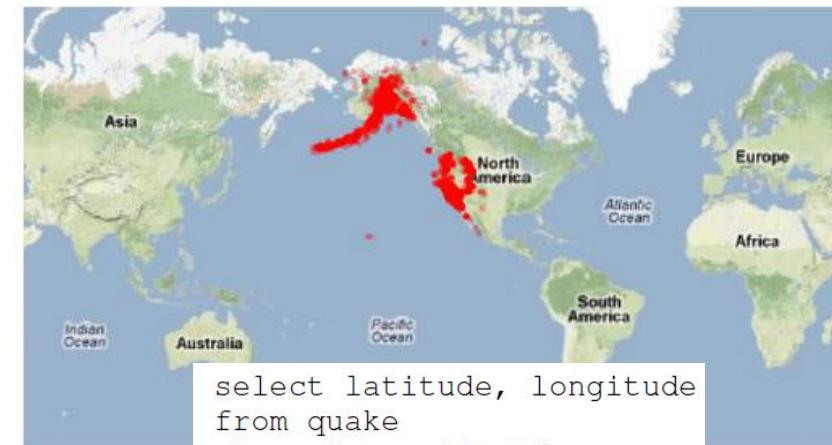
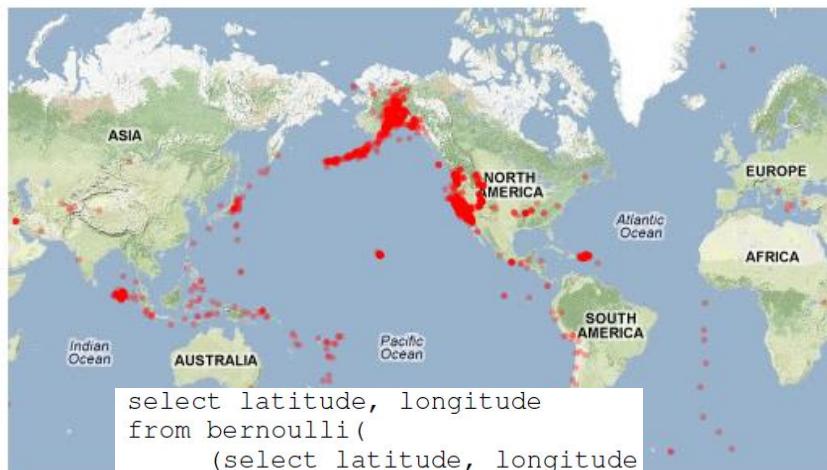
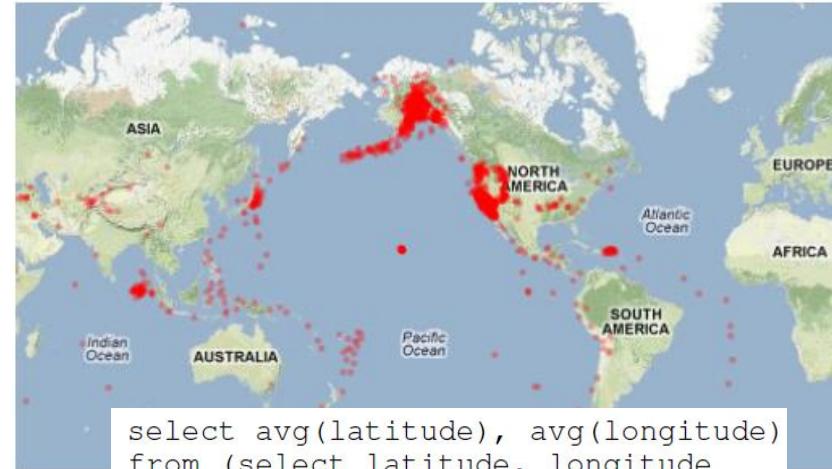
ScalaR: resolution reduction

Resolution Reduction Techniques:

- **Aggregation:** group the data into n sub-matrices and return summaries over the sub-matrices
- **Sampling:** given a probability value p, return roughly that fraction of data as the result, where $p * |data| = n$
- **Filtering:** given a set of filters over the data, return the elements that pass these filters

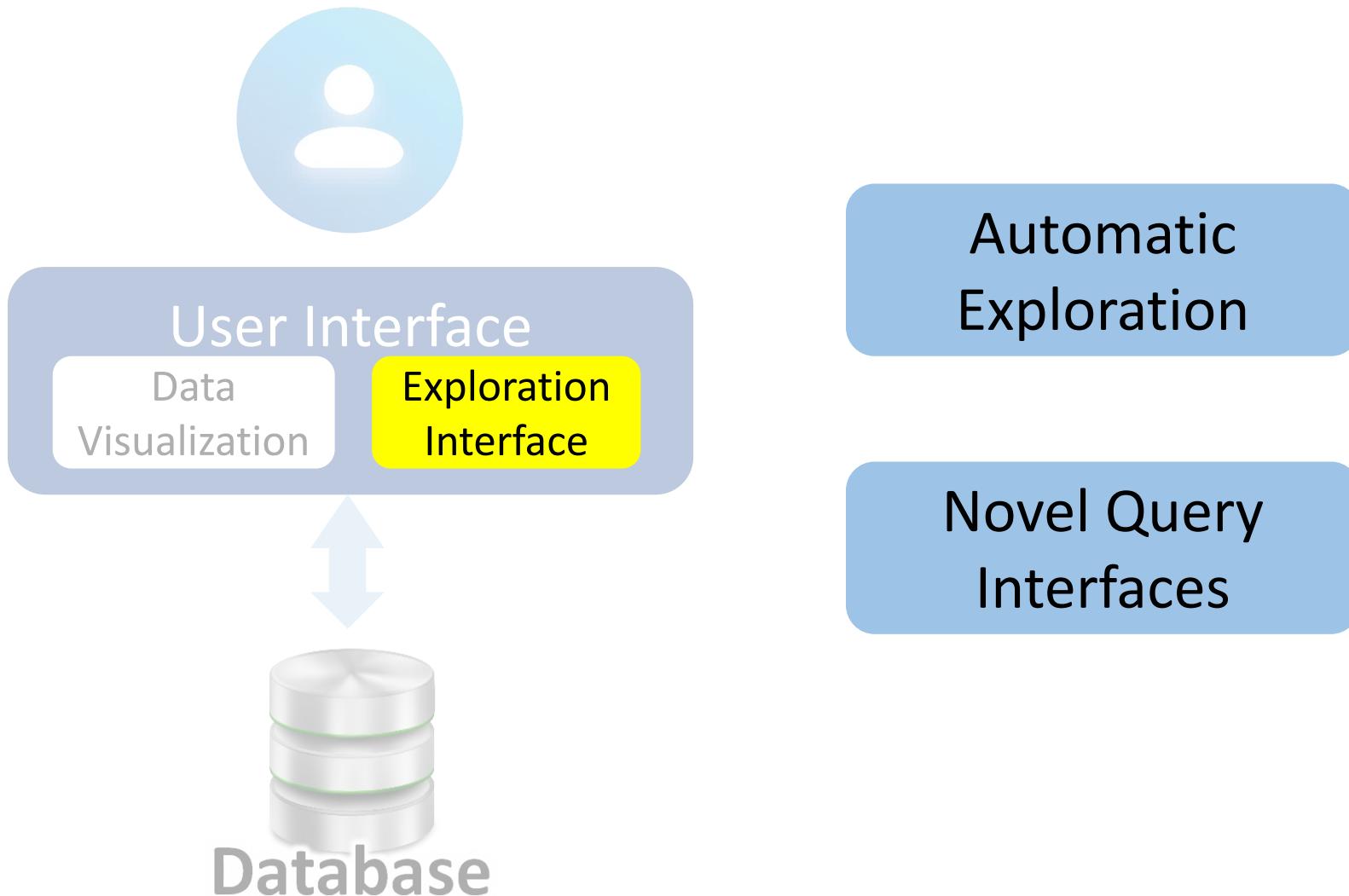


6381*6543 (a) Original query





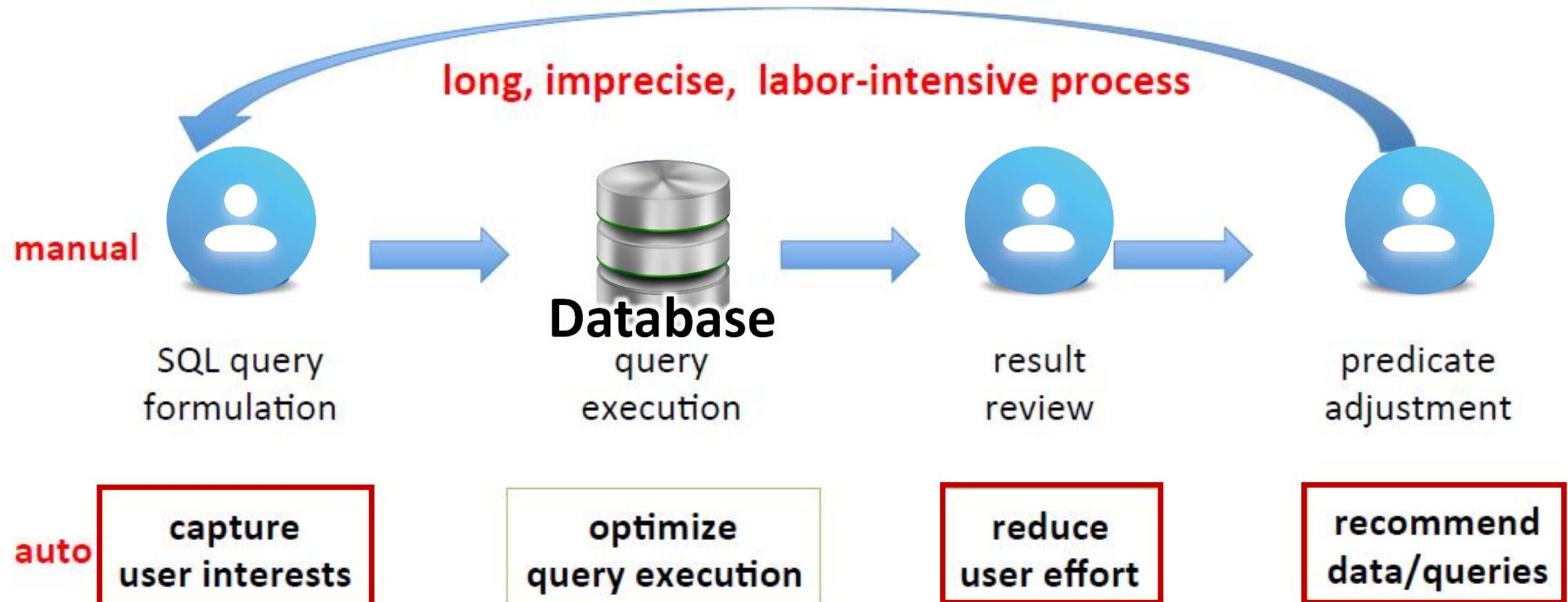
IDE: Exploration Interface



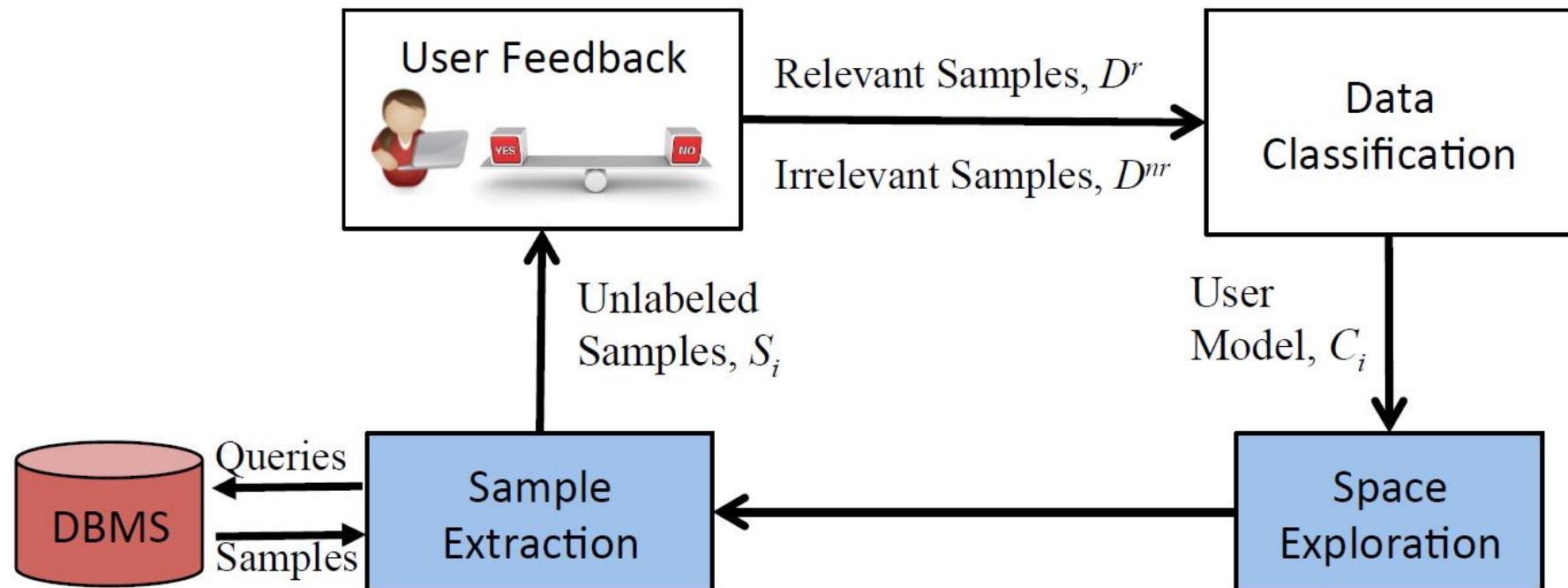


Automatic Exploration

Manual vs Automatic Exploration



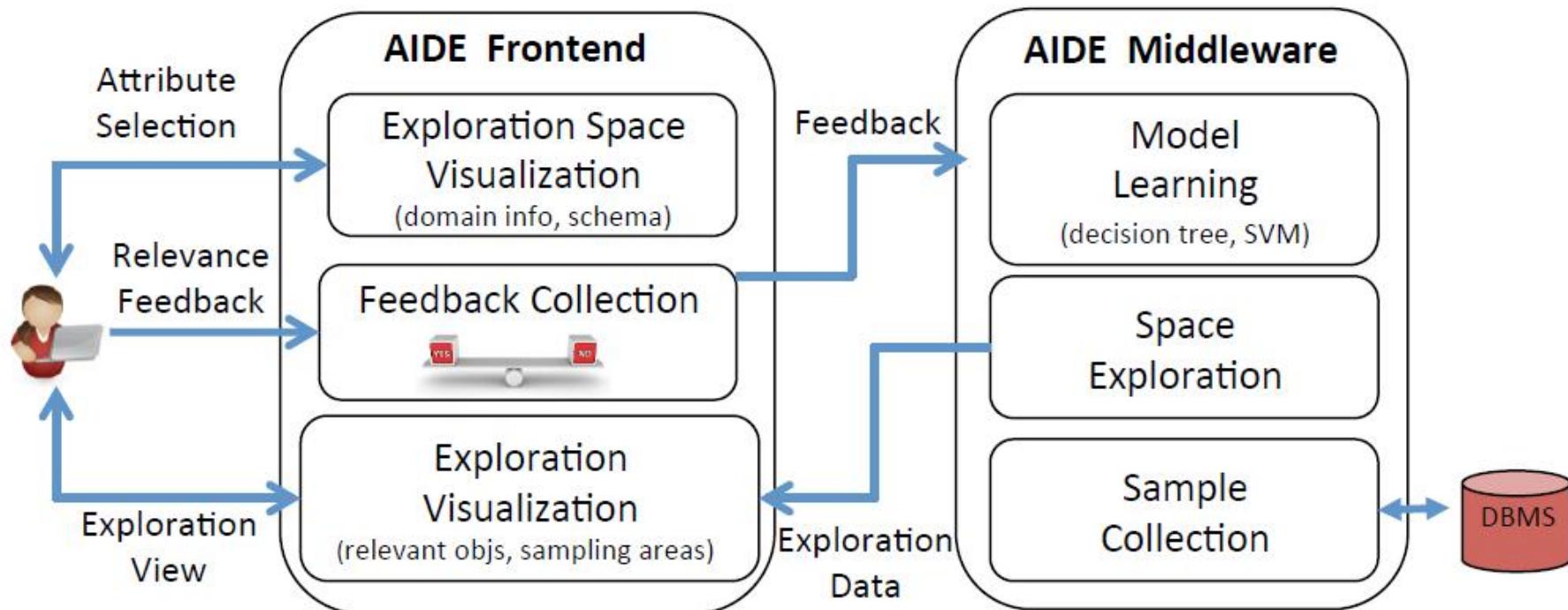
Explore by Example



AIDE: an automatic user navigation system for interactive data exploration. Diao et al. VLDB 2015.

AIDE: an active learning-based approach for interactive data exploration. K. Dimitriadou, O. Papaemmanouil, Y. Diao. TKDE 2016.

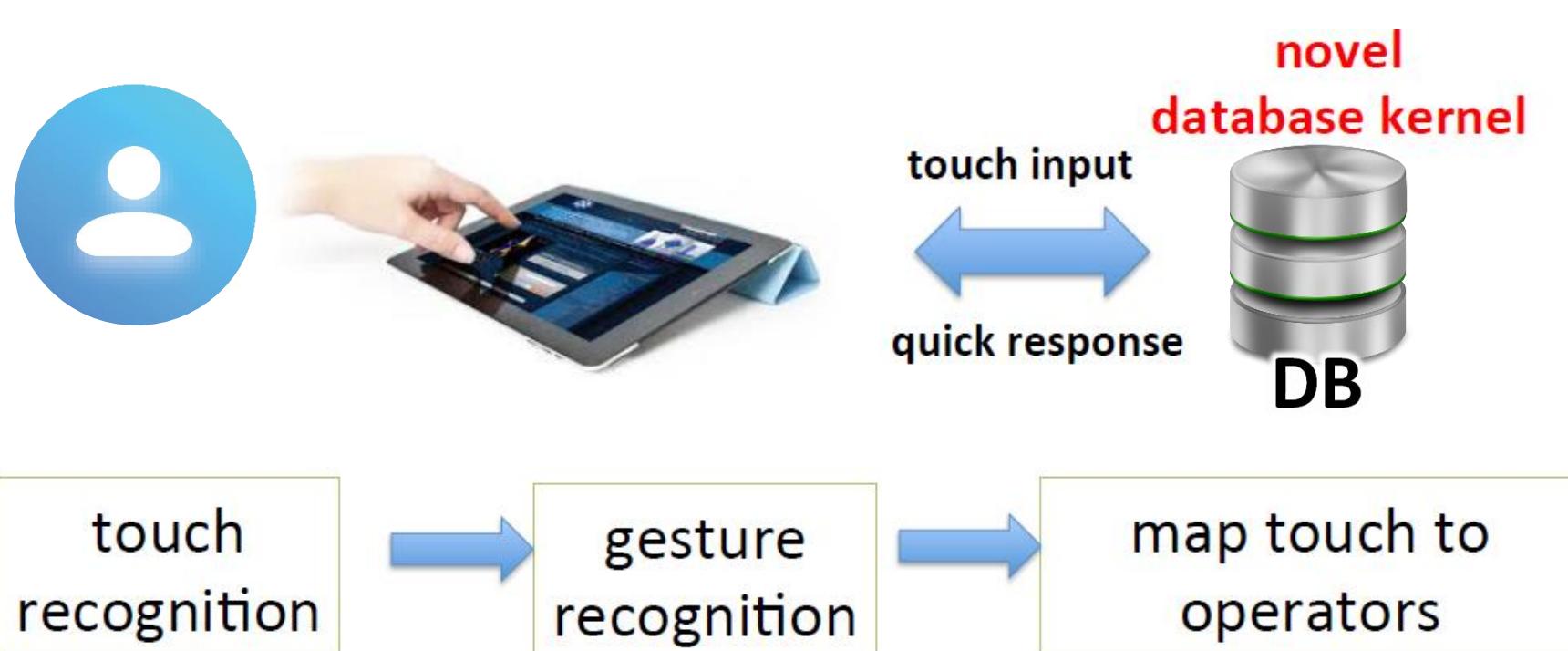
Explore by Example



Query Interfaces



No-keyboard Interfaces



No-keyboard Interfaces

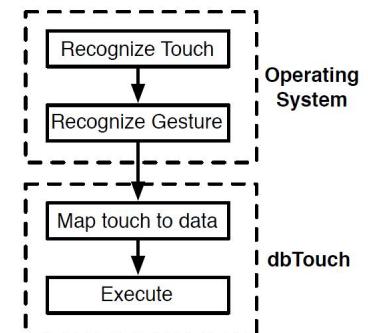
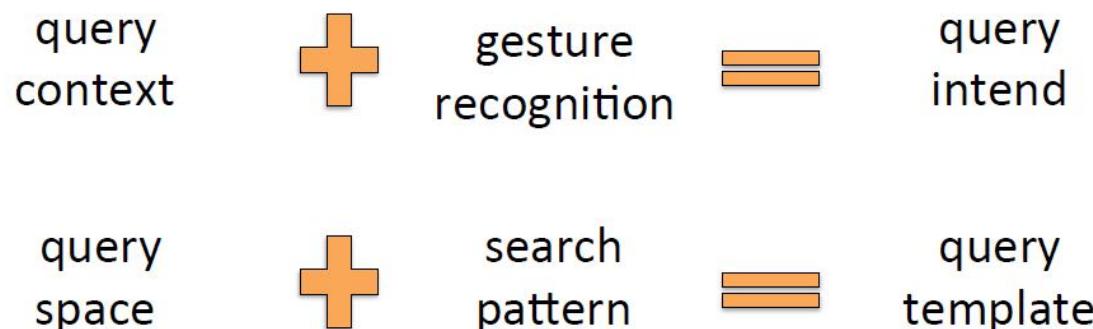
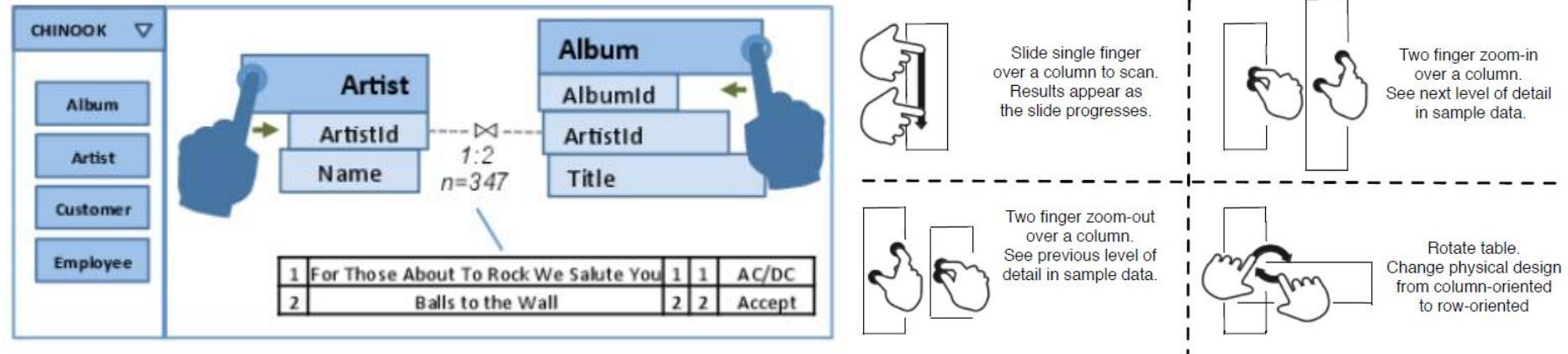
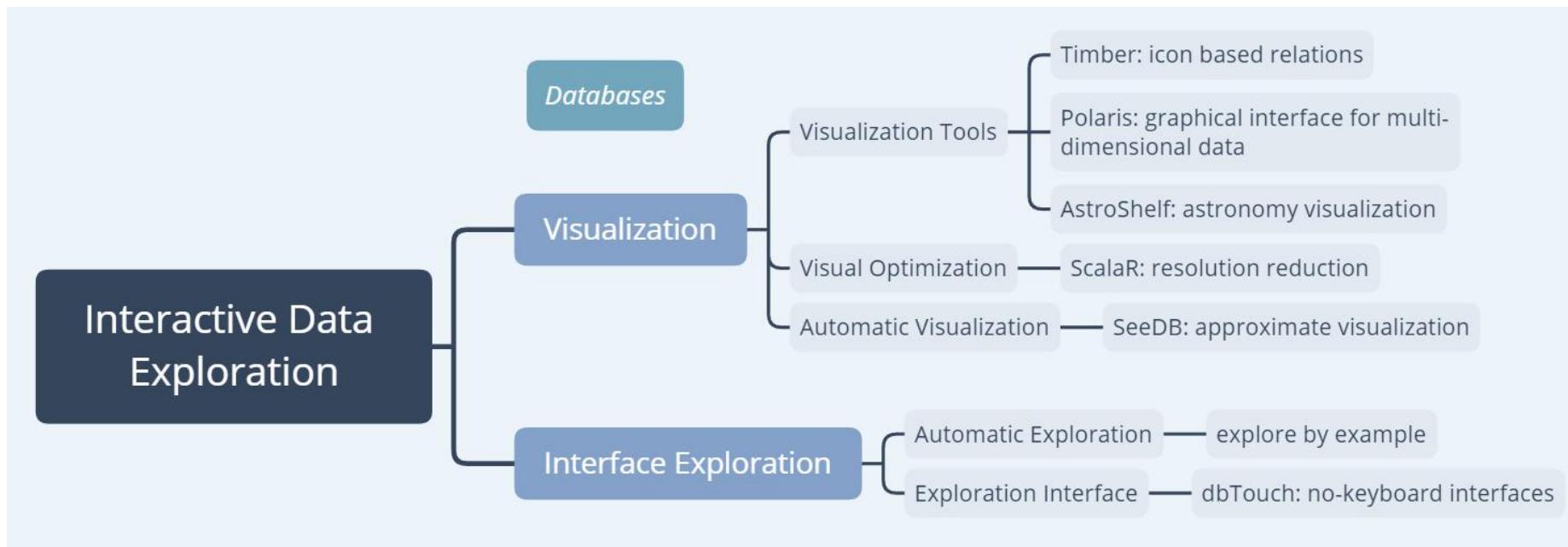


Figure 3: dbTouch system layers.

Summary



- Big data is significant for artificial intelligence and human intelligence
- Project-driven course
- Interactive data exploration systems



Thank You

QQ群: 1057859502

