

# 大数据综合处理实验 课程简介

鸣谢：本课程得到Google与Intel公司  
中国大学合作部精品课程计划资助

南京大学计算机科学与技术系

教师：黄宜华，顾荣

# 课程简介

## 教学内容简介

本课程将系统介绍大数据并行处理技术和编程方法。课程首先介绍并行计算技术的基本概念、原理、方法和技术，在此基础上，介绍基于集群的大数据并行处理技术原理和方法，重点介绍Hadoop MapReduce并行计算集群的构架、用于大数据存储和计算的分布式文件系统、以及基于MapReduce集群的大数据并行处理技术和编程方法，MapReduce并行化算法设计技术、并行化算法应用研究案例。

## 教学目标

课程的主要目标是通过介绍多处理器并行处理技术、以及基于集群的大数据并行处理技术和MapReduce并行编程模型和方法，要求学生理解和掌握并行处理技术的基本概念、原理和构架、以及基于集群的大规模海量数据并行处理与编程技术方法，并能够用MapReduce实际设计和编写具体的大数据处理应用问题的算法和程序。

## 选课要求

具有Java程序设计能力、Linux系统操作使用能力、以及基础性机器学习与数据挖掘算法知识；除课堂听课外需要完成编程实验，还要求在学期结束时完成一个综合性课程设计

# 选修本课程的重要性

## 并行处理成为计算技术的重大发展趋势

- 单处理器性能提升达到极限, 多核/多处理并行计算成为计算技术发展必然趋势
- 并行计算技术将渗透到每个计算应用领域
- 并行计算技术将影响传统计算技术的各个层面, 与传统计算技术相结合产生很多新的研究热点和课题

## IT行业和应用已进入“大数据(Big Data)”和“数据为王”的时代

IT行业应用规模急剧扩大, 出现越来越多的超大规模数据处理应用需求, 传统系统难以提供足够的存储和计算资源进行处理

- 2008年国际著名的《Nature》杂志出版一期专刊专门讨论未来大数据(Big Data)处理相关的技术问题和挑战
- 世界权威的IT数据分析公司IDC: 全世界的数据量2009年为800EB, 到2020年将增长44倍, 达到35ZB(35, 000EB)
- 未来的IT行业中, 价值在于数据, 未来是“数据为王”的时代

# 选修本课程的重要性

计算机专业人员面临挑战，市场迫切需要相应的专业技术人才

- 并行计算技术将从硬件到软件全面影响传统计算技术的各个层面，新的技术发展挑战和需求迫使我们软件开发和程序设计人员必须尽快掌握并行计算技术
- 20-30年前程序设计技术最大的革命是面向对象技术，而下一个程序设计技术的革命将是并行程序设计技术
- 今天绝大多数程序员不懂并行设计技术，就像20年前绝大多数程序员不懂面向对象技术一样
- 目前国内外的知名IT企业迫切需要大量掌握大规模数据并行处理技术的人才

# 课程内容

## Ch.1 并行计算技术简介

简要介绍并行计算技术的概况，基本分类，主要技术问题，MPI并行程序设计，大规模并行数据处理技术

## Ch.2 MapReduce简介

简要介绍MapReduce技术的由来，基本构思，编程模型，主要设计思想和技术特征，基本应用

## Ch.3 Google 和Hadoop MapReduce的基本构架

介绍Google MapReduce并行计算框架的基本结构、工作原理，Google分布式文件系统GFS的基本构架与工作原理，Google结构化数据管理系统BigTable的基本结构与工作原理

介绍开源大数据处理系统Hadoop 的基本组成结构和工作原理，MapReduce基本框架和工作原理，HDFS基本组成及工作原理，并介绍HDFS的基本编程

# 课程内容

## Ch.4 Hadoop系统安装运行与程序开发

介绍单机和集群Hadoop系统安装方法和步骤，以及程序开发环境与开发过程

### 实验1: Hadoop系统安装与WordCount词频统计编程实验

## Ch.5 MapReduce算法设计

介绍排序算法、文档倒排索引、文档共现算法、专利文献数据分析应用

### 实验2: 搜索引擎文档倒排索引编程实验

## Ch.6 Hadoop HBase与Hive原理与编程技术

介绍Hadoop 分布式数据管理系统HBase工作原理及其编程技术；介绍Hadoop数据仓库 Hive基本结构、工作原理及其编程技术

### 实验3: Hadoop HBase和Hive编程实验（待定）

# 课程内容

## Ch.7 高级MapReduce编程技术

介绍复杂I/O数据表示、用复合键值对完成特殊处理、程序员定制的I/O格式、Partitioner、Combiner，基于迭代的MapReduce求解方法、数据相关MapReduce任务计算、链式MapReduce计算、多数据源连接、访问关系数据库等高级技术

## Ch.8 基于MapReduce的搜索引擎算法

介绍网页排名算法PageRank，搜索引擎文档倒排索引算法，以及全文检索系统的设计实现

**实验4：Wikipedia网页PageRank实验（待定）**

## Ch.9 基于MapReduce的数据挖掘基础算法

介绍机器学习和数据挖掘中的聚类算法、分类算法、频繁项集挖掘等算法的MapReduce并行化设计技术方法

**实验5：待定**

# 课程内容

## Ch.10 Spark系统和编程技术介绍

介绍基于内存计算的Spark系统及其基本编程技术

## 课程设计大作业

自选或由老师指定具有一定难度和工作量的题目，完成一个综合性大数据课程设计



# 课时安排

2020年春季学期

**学期中：** 课堂讲授，课程实验，复习，期末考试  
每周2课时，共计18次（32课时）

**暑假7月份：** 分组完成综合课程设计

# 考核方法

**期末考试**  
笔试，占50%

**课程实验**  
5次，共计占25%

**课程设计**  
占25%

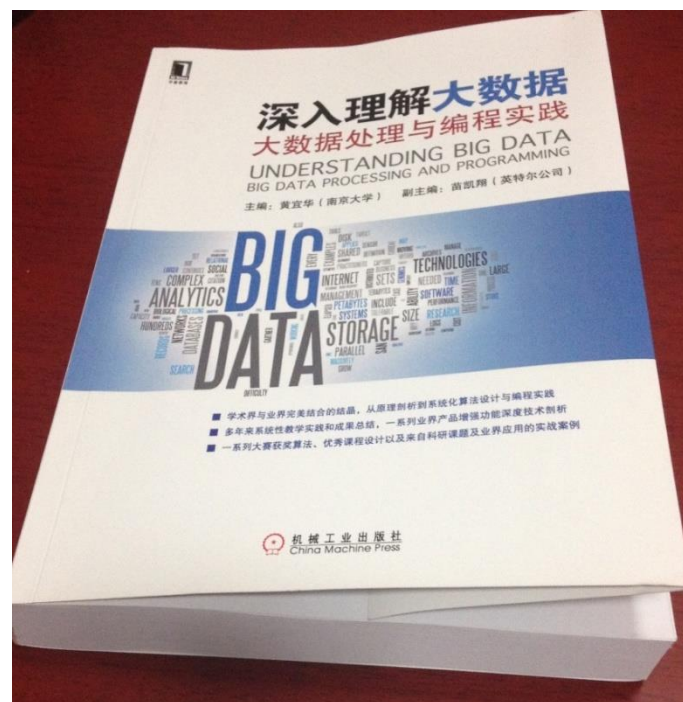
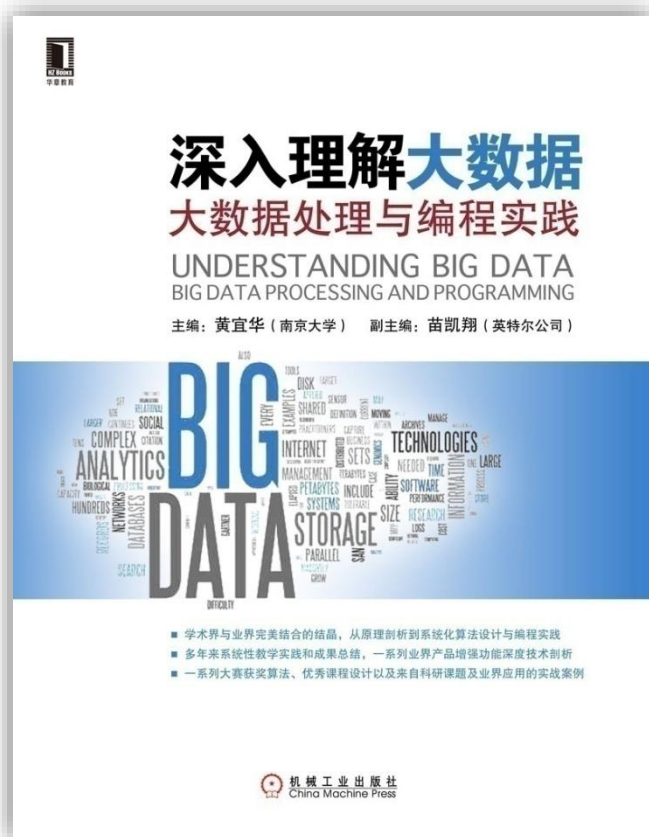
# 课件与参考书目、文献

课件参考资料下载：计算机系本科教学服务平台

教材（安装配置，编程API部分版本较老，最新的参考相关软件官网）：

《深入理解大数据—大数据编程技术与实践》

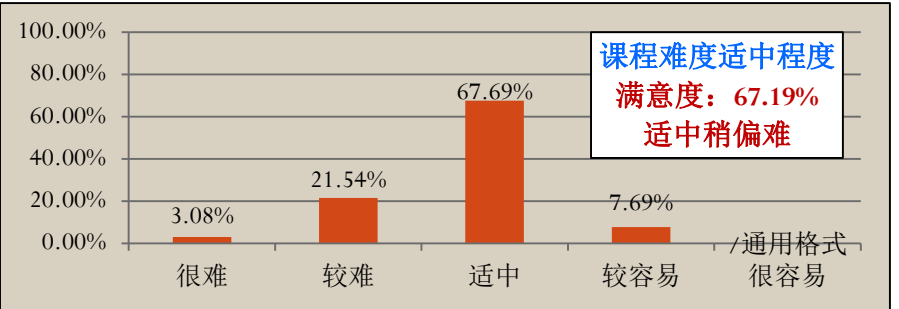
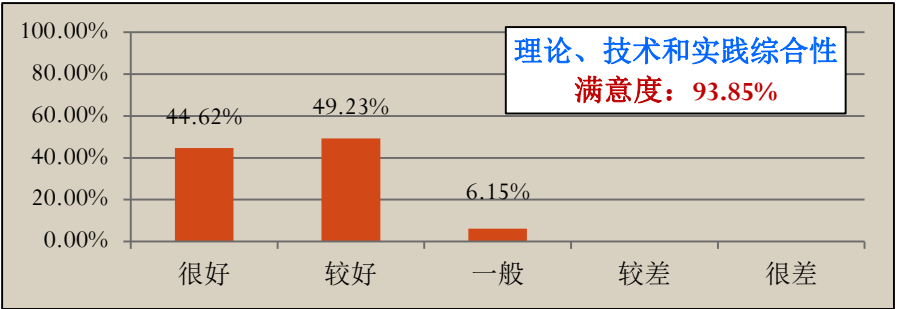
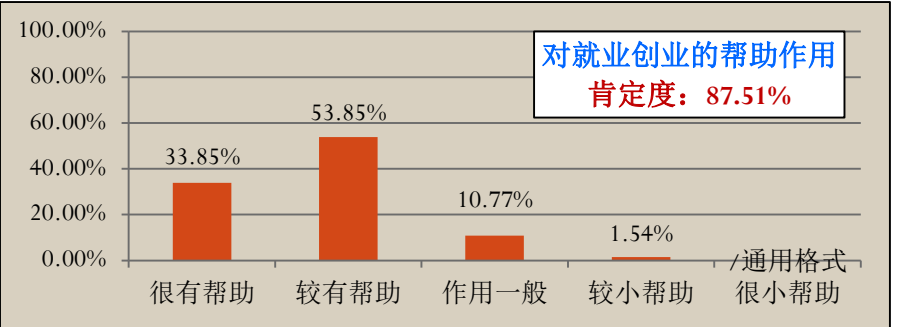
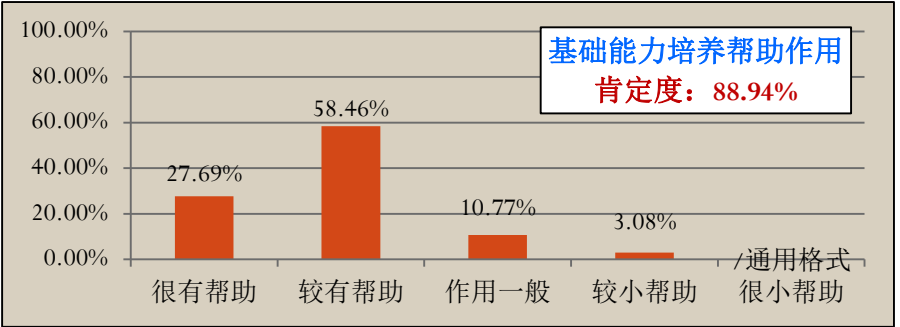
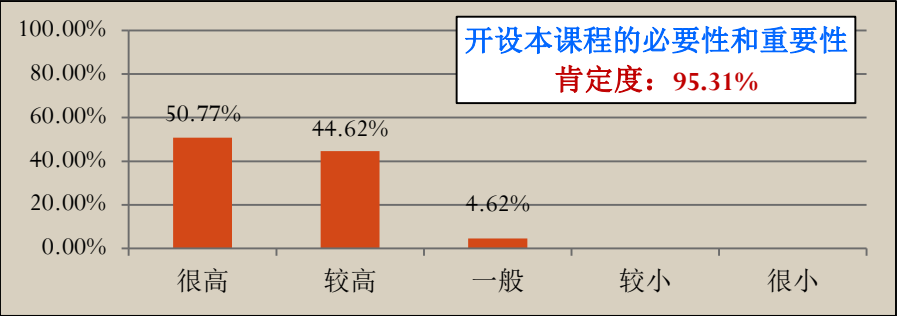
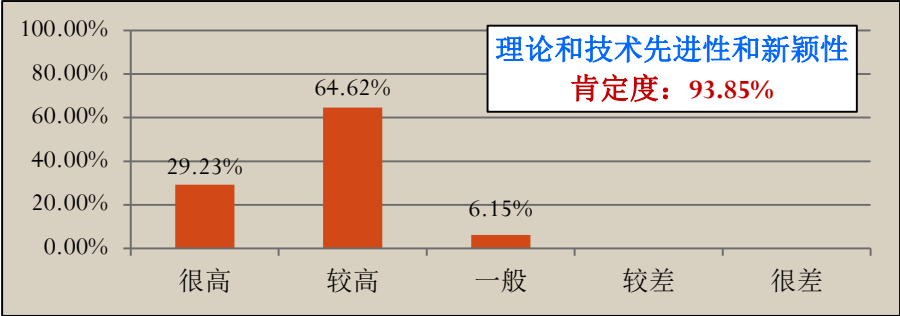
黄宜华，苗凯翔主编，机械工业出版社，2014



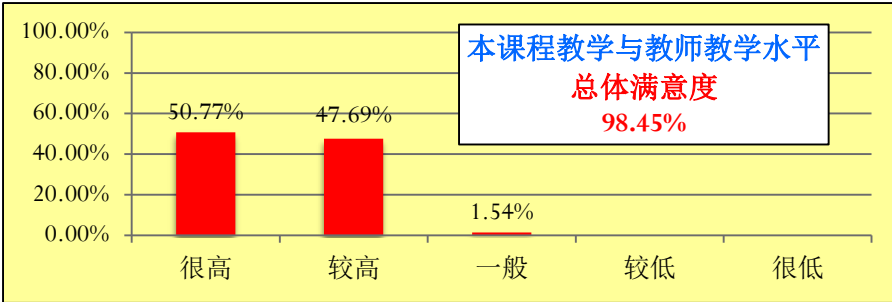
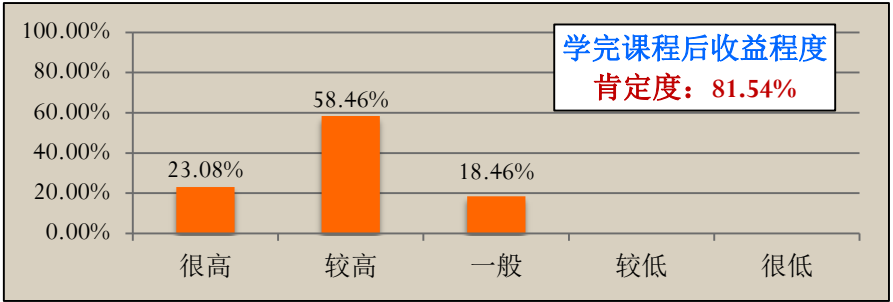
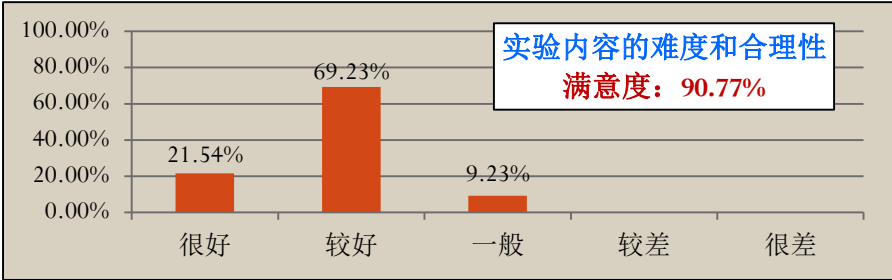
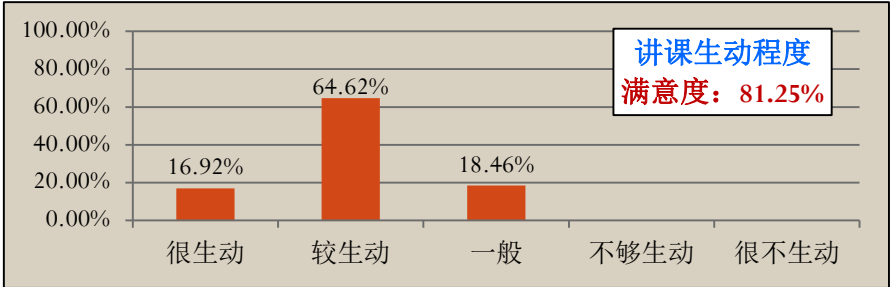
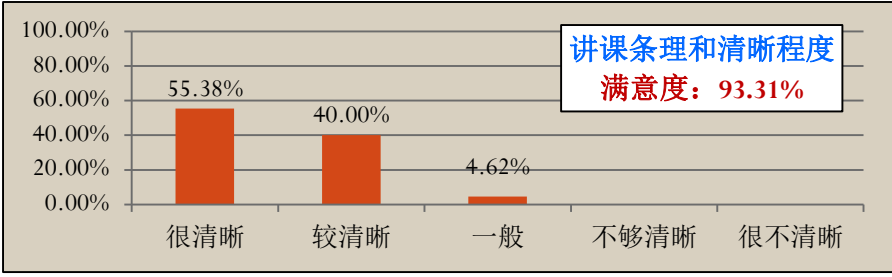
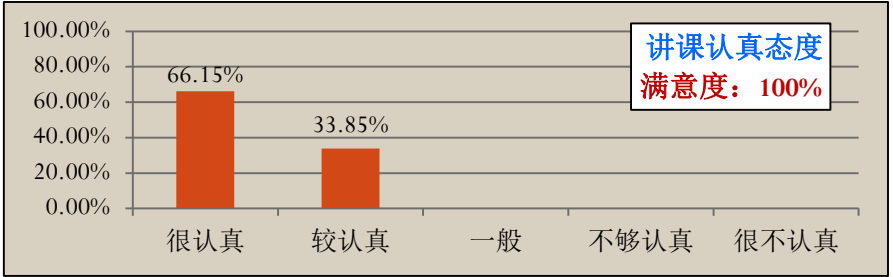
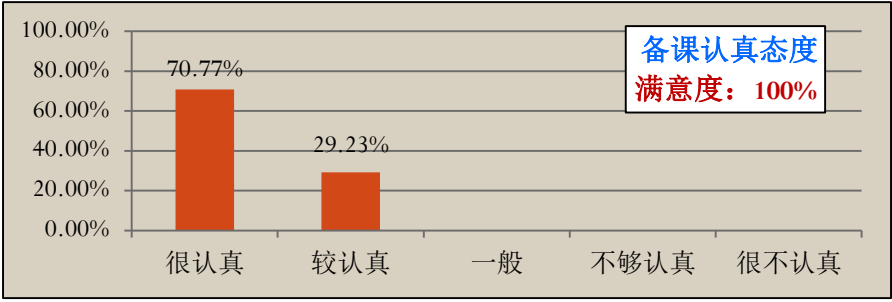
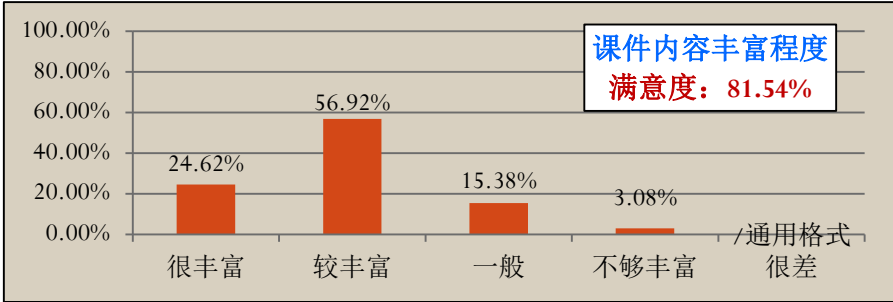
# 课程开课情况

- 2011-2018学年已经为研究生开设了8个学期, 每周2学时, 研究生和本科生, 本系选修人数900多人 (研究生为主)
- 2016年首次为本科生开设大数据综合实验课程, 120人选修, 课程情况良好, 绝大多数同学一次通过, 仅有6位同学因笔试成绩不好或因课程设计综合实验未做, 后补考通过
- 安排5次又简到难的课程实验, 从基本工具平台使用到编程实验
- 课程结束后, 要求分小组完成一个具有一定难度的课程项目设计, 每年都会出现一批相当出色的课程设计项目。

# 课程教学评估



# 课程教学评估



# 教学效果

## 第一届“中国云/移动互联网创新大奖赛”



本课程开设后，我系机器学习与数据挖掘研究所和云计算与大数据并行计算课题组学习了MapReduce技术的同学组织了4支研究生代表队在“中国云产业联盟”组织的首届“中国云·移动互联网创新大奖赛”中参赛并荣获9项优胜奖（一等奖2项，二等奖4项，三等奖3项）和4项优秀领队奖，并获得大赛奖金20万元！占据大赛全部30个奖项中的9项，4道大数据赛题全部17个奖项中的8项！共有来自全国高校的230多支队伍参赛。



# 教学效果

## 第一届“中国云/移动互联网创新大奖赛”



技术类赛题 1: 调色板搜图—在百万图片中搜索与指定调色板相近的图片

技术类赛题 2: 多快好省的速递员 — 动态路况环境下的物流规划

技术类赛题 3: 你不知道我知道 — 互联网问答系统用户行为分析

技术类赛题 4: 难舍难分 — 大规模搜索关键字（短文本）分类

技术类赛题 5: 麻雀级云数据中心 — 规定时间内在小规模硬件环境上部署大量虚拟机

创意类竞赛说明: 创意类赛题没有具体的问题约束。



# 教学效果

## 第一届“中国云/移动互联网创新大奖赛”



我系4支研究生代表队荣获9项优胜奖和4项优秀领队奖，获得奖金20万



# 教学效果

## 第一届“中国云/移动互联网创新大奖赛”

“中国云·移动互联网创新大奖赛”是由“中国云产业联盟”和百度、阿里巴巴、腾讯、用友等国内著名企业和北航、北大等著名高校于2012年5月联合发起组织的第一届全国云计算和互联网创新技术大赛。这次大赛由北航的怀俊鹏院士与中国云产业联盟联合倡议并发起，来自国内多所著名高校和著名企业的十多位专家学者共同参与，是目前为止国内规模和影响最大、级别最高的云计算和互联网创新技术大赛。颁奖仪式上，中国科学院院士怀俊鹏教授、微软集团副总裁陆奇博士到会做了关于云计算的主题报告，百度总裁李彦宏、宽带资本董事长田溯宁、用友软件董事长兼总裁王文京、中国联通总裁陆益民等嘉宾也到会并做了云计算主题对话

### » 参会嘉宾（排名不分先后）



北航校长  
怀俊鹏

#### 大数据时代面临三大挑战

1. 软件 and 数据处理能力。  
2. 资源和共享管理的挑战。  
3. 数据处理的可信能力。...[全文]



微软集团副总裁  
陆奇

#### 下个时代是智能交互的时代

下个时代是智能交互的时代，通过深度机器学习，机器可以理解人类的语言，机器还可以理解人类的手势语言，机器可以像人一样观察世界。...[全文]



宽带资本董事长  
田溯宁

云产业联盟，中国云的理想，是需要大家的努力，每个大赛的团队，每个创意的思想，都是播下的种子，我们希望这个种子能够随着中国的经济发展，随着我们每个企业的发展，能够茁壮成长，能够成为参天大树，能够建立中国云计算，中国大数据的生态系统。[全文]



百度CEO  
李彦宏

在百度的眼中，运营商是一个巨大的未开采的金矿。运营商掌握的数据是我们梦寐以求的数据。百度积累的云计算的能力，运营商过去积累的系统的数，用户很好的结合起来的话，可以产生很多新的创新。...[全文]



中国联通总裁  
陆益民

未来运营商的出路在哪儿？一定也是创新。未来的创新可能也是在于商业模式的创新，技术的创新。云计算是给我们创新的一个很重要的方向。云战略是作为我们未来战略很重要的方向。中国联通在这个方面。...[全文]



用友软件董事长兼总裁  
王文京

从企业市场来讲，无论是大中型企业，还是小微企业，都有巨大的机会。加上整个中国市场规模的有利条件，我特别赞同刚才陆博士讲的，这是一个中国云的时代，这样的历史机会已经到来。...[全文]



龙湖地产董事长  
吴亚军

我们公司的使命就是除了提供优质的产品和服务之外，我们期待影响他人的行为。我们最想研究人的行为，但是现在很难采集这样的数据。第二，我们拿到这些数据，在加工和分析的过程中，其实我们也有诸多的困难。...[全文]



云联盟秘书长  
姜广智

为鼓励首届iCom大赛获奖选手，持续开展技术创新活动，云联盟企业庄严承诺，为所有获得首届大赛奖励的选手，提供直接就业机会，暨大赛获奖选手获奖元件和本承诺元件直接进入企业进行就业，获得同等条件下优先录取的条件。...[全文]

# 教学效果

## “全国高校云计算创新应用大赛” 大数据技能赛

- 2015-2018年，荣获教育部主办的“全国高校云计算创新应用大赛”大数据技能赛冠军，实现该项赛事全国四连冠！



## 本科课程

今年是本系第四次作为本科生准出核心平台选修课开设本课程

希望大家积极配合以便为大家开好本课程，圆满完成本课程的教学工作，达到预定的教学目标！

**谢谢大家！**