



哈尔滨工业大学

海量数据计算研究中心

Massive Data Computing Lab @ HIT

大数据算法

第一讲 大数据算法概述

哈尔滨工业大学

王宏志

wangzh@hit.edu.cn

本讲内容

- 1.1 大数据的定义与特点
- 1.2 大数据算法
- 1.3 大数据算法设计与分析

什么是大数据？

- 至今没有公认的定义

- 定义1 (Kusnetzky, Dan. What is "Big Data?")

所涉及的数据量规模巨大到无法通过人工，在合理时间内达到截取、管理、处理、并整理成为人类所能解读的信息

- 定义2 (维克托·迈尔-舍恩伯格、肯尼斯·库克耶. “大数据时代”)

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

- 定义3 (“大数据” (Big data) 研究机构Gartner)

“大数据”是需要新处理模式才能具有更强的决策力、洞察发现力和流程优化能力的海量、高增长率和多样化的信息资产。

“大” 有多大

1Byte = 8 bit

1 KB = 1,024 Bytes

1 MB = 1,024 KB = 1,048,576 Bytes

1 GB = 1,024 MB = 1,048,576 KB

1 TB = 1,024 GB = 1,048,576 MB

1 PB = 1,024 TB = 1,048,576 GB

1 EB = 1,024 PB = 1,048,576 TB

1 ZB = 1,024 EB = 1,048,576 PB

1 YB = 1,024 ZB = 1,048,576 EB

1 BB = 1,024 YB = 1,048,576 ZB

1 NB = 1,024 BB = 1,048,576 YB

1 DB = 1,024 NB = 1,048,576 BB

仟 1,000

萬 10^4

億 10^8

兆 10^{12}

京 10^{16}

垓 10^{20}

秭 10^{24}

穰 10^{28}

沟 10^{32}

涧 10^{36}

正 10^{40}

载 10^{44}

俱胝 10^7

阿庾多 10^{14}

那由他 10^{28}

频波罗 10^{56}

矜羯罗 10^{112}

阿伽罗 10^{224}

最胜 10^{448}

摩婆罗 10^{896}

阿婆罗 10^{1792}

多婆罗 10^{3584}

界分 10^{7168}

普摩 10^{14336}

祢摩 10^{28672}

阿婆铃 10^{57344}

弥伽婆 10^{114688}

毗攞伽 10^{229376}

毗伽婆 10^{458752}

僧羯逻摩 10^{917504}

毗萨罗 $10^{1835008}$

毗赡婆 $10^{3670016}$

毗盛伽 $10^{7340032}$

毗素陀 $10^{14680064}$

毗婆诃 $10^{29360128}$

毗薄底 $10^{58720256}$

毗佉担 $10^{117440512}$

称量 $10^{234881024}$

一持 $10^{469762048}$

异路 $10^{939524096}$

颠倒 $10^{1879048192}$

三末耶 $10^{3758096384}$

毗睹罗 $10^{7516192768}$

奚婆罗 $10^{15032385536}$

伺察 $10^{30064771072}$

周广 $10^{60129542144}$

高出 $10^{120259084288}$

最妙 $10^{240518168576}$

泥罗婆 $10^{481036337152}$

诃理婆 $10^{962072674304}$

一动 $10^{1924145348608}$

诃理蒲 $10^{3848290697216}$

诃理三 $10^{7696581394432}$

奚鲁伽 $10^{15393162788864}$

达攞步陀 $10^{30786325577728}$

诃鲁那 $10^{61572651155456}$

摩鲁陀 $10^{123145302310912}$

忉慕陀 $10^{246290604621824}$

瑩攞陀 $10^{492581209243648}$

摩鲁摩 $10^{985162418487296}$

“大” 有多大 (续)

调伏 10¹⁹⁷⁰³²⁴⁸³⁶⁹⁷⁴⁵⁹²

离憍慢 10³⁹⁴⁰⁶⁴⁹⁶⁷³⁹⁴⁹¹⁸⁴

不动 10⁷⁸⁸¹²⁹⁹³⁴⁷⁸⁹⁸³⁶⁸

极量 10¹⁵⁷⁶²⁵⁹⁸⁶⁹⁵⁷⁹⁶⁷³⁶

阿么怛罗 10³¹⁵²⁵¹⁹⁷³⁹¹⁵⁹³⁴⁷²

勃么怛罗 10⁶³⁰⁵⁰³⁹⁴⁷⁸³¹⁸⁶⁹⁴⁴

伽么怛罗 10¹²⁶¹⁰⁰⁷⁸⁹⁵⁶⁶³⁷³⁸⁸⁸

那么怛罗 10²⁵²²⁰¹⁵⁷⁹¹³²⁷⁴⁷⁷⁷⁶

奚么怛罗 10⁵⁰⁴⁴⁰³¹⁵⁸²⁶⁵⁴⁹⁵⁵⁵²

鞞么怛罗 10¹⁰⁰⁸⁸⁰⁶³¹⁶⁵³⁰⁹⁹¹¹⁰⁴

鉢罗么怛罗 10²⁰¹⁷⁶¹²⁶³³⁰⁶¹⁹⁸²²⁰⁸

尸婆么怛罗 10⁴⁰³⁵²²⁵²⁶⁶¹²³⁹⁶⁴⁴¹⁶

翳罗 10⁸⁰⁷⁰⁴⁵⁰⁵³²²⁴⁷⁹²⁸⁸³²

薜罗 10¹⁶¹⁴⁰⁹⁰¹⁰⁶⁴⁴⁹⁵⁸⁵⁷⁶⁶⁴

谛罗 10³²²⁸¹⁸⁰²¹²⁸⁹⁹¹⁷¹⁵³²⁸

偈罗 10⁶⁴⁵⁶³⁶⁰⁴²⁵⁷⁹⁸³⁴³⁰⁶⁵⁶

宰步罗 10¹²⁹¹²⁷²⁰⁸⁵¹⁵⁹⁶⁶⁸⁶¹³¹²

泥罗 10²⁵⁸²⁵⁴⁴¹⁷⁰³¹⁹³³⁷²²⁶²⁴

计罗 10⁵¹⁶⁵⁰⁸⁸³⁴⁰⁶³⁸⁶⁷⁴⁴⁵²⁴⁸

细罗 10¹⁰³³⁰¹⁷⁶⁶⁸¹²⁷⁷³⁴⁸⁹⁰⁴⁹⁶

睥罗 10²⁰⁶⁶⁰³⁵³³⁶²⁵⁵⁴⁶⁹⁷⁸⁰⁹⁹²

谜罗 10⁴¹³²⁰⁷⁰⁶⁷²⁵¹⁰⁹³⁹⁵⁶¹⁹⁸⁴

娑攞荼 10⁸²⁶⁴¹⁴¹³⁴⁵⁰²¹⁸⁷⁹¹²³⁹⁶⁸

谜鲁陀 10¹⁶⁵²⁸²⁸²⁶⁹⁰⁰⁴³⁷⁵⁸²⁴⁷⁹³⁶

契鲁陀 10³³⁰⁵⁶⁵⁶⁵³⁸⁰⁰⁸⁷⁵¹⁶⁴⁹⁵⁸⁷²

摩睹罗 10⁶⁶¹¹³¹³⁰⁷⁶⁰¹⁷⁵⁰³²⁹⁹¹⁷⁴⁴

娑母罗 10¹³²²²⁶²⁶¹⁵²⁰³⁵⁰⁰⁶⁵⁹⁸³⁴⁸⁸

阿野娑 10²⁶⁴⁴⁵²⁵²³⁰⁴⁰⁷⁰⁰¹³¹⁹⁶⁶⁹⁷⁶

迦么罗 10⁵²⁸⁹⁰⁵⁰⁴⁶⁰⁸¹⁴⁰⁰²⁶³⁹³³⁹⁵²

摩伽婆 10¹⁰⁵⁷⁸¹⁰⁰⁹²¹⁶²⁸⁰⁰⁵²⁷⁸⁶⁷⁹⁰⁴

阿怛罗 10²¹¹⁵⁶²⁰¹⁸⁴³²⁵⁶⁰¹⁰⁵⁵⁷³⁵⁸⁰⁸

酰鲁耶 10⁴²³¹²⁴⁰³⁶⁸⁶⁵¹²⁰²¹¹¹⁴⁷¹⁶¹⁶

薜鲁婆 10⁸⁴⁶²⁴⁸⁰⁷³⁷³⁰²⁴⁰⁴²²²⁹⁴³²³²

羯罗波 10¹⁶⁹²⁴⁹⁶¹⁴⁷⁴⁶⁰⁴⁸⁰⁸⁴⁴⁵⁸⁸⁶⁴⁶⁴

诃婆婆 10³³⁸⁴⁹⁹²²⁹⁴⁹²⁰⁹⁶¹⁶⁸⁹¹⁷⁷²⁹²⁸

毗婆罗 10⁶⁷⁶⁹⁹⁸⁴⁵⁸⁹⁸⁴¹⁹²³³⁷⁸³⁵⁴⁵⁸⁵⁶

那婆罗 10¹³⁵³⁹⁹⁶⁹¹⁷⁹⁶⁸³⁸⁴⁶⁷⁵⁶⁷⁰⁹¹⁷¹²

摩攞罗 10²⁷⁰⁷⁹⁹³⁸³⁵⁹³⁶⁷⁶⁹³⁵¹³⁴¹⁸³⁴²⁴

娑婆罗 10⁵⁴¹⁵⁹⁸⁷⁶⁷¹⁸⁷³⁵³⁸⁷⁰²⁶⁸³⁶⁶⁸⁴⁸

迷攞普 10¹⁰⁸³¹⁹⁷⁵³⁴³⁷⁴⁷⁰⁷⁷⁴⁰⁵³⁶⁷³³⁶⁹⁶

者么罗 10²¹⁶⁶³⁹⁵⁰⁶⁸⁷⁴⁹⁴¹⁵⁴⁸¹⁰⁷³⁴⁶⁷³⁹²

馱么罗 10⁴³³²⁷⁹⁰¹³⁷⁴⁹⁸⁸³⁰⁹⁶²¹⁴⁶⁹³⁴⁷⁸⁴

鉢攞么陀 10⁸⁶⁶⁵⁵⁸⁰²⁷⁴⁹⁹⁷⁶⁶¹⁹²⁴²⁹³⁸⁶⁹⁵⁶⁸

毗迦摩 10¹⁷³³¹¹⁶⁰⁵⁴⁹⁹⁹⁵³²³⁸⁴⁸⁵⁸⁷⁷³⁹¹³⁶

乌波跋多 10³⁴⁶⁶²³²¹⁰⁹⁹⁹⁹⁰⁶⁴⁷⁶⁹⁷¹⁷⁵⁴⁷⁸²⁷²

演说 10⁶⁹³²⁴⁶⁴²¹⁹⁹⁹⁸¹²⁹⁵³⁹⁴³⁵⁰⁹⁵⁶⁵⁴⁴

无尽 10¹³⁸⁶⁴⁹²⁸⁴³⁹⁹⁹⁶²⁵⁹⁰⁷⁸⁸⁷⁰¹⁹¹³⁰⁸⁸

出生 10²⁷⁷²⁹⁸⁵⁶⁸⁷⁹⁹⁹²⁵¹⁸¹⁵⁷⁷⁴⁰³⁸²⁶¹⁷⁶

“大” 有多大 (续)

无我 $10^{554597137599850363154807652352}$

阿畔多 $10^{1109194275199700726309615304704}$

青莲华 $10^{2218388550399401452619230609408}$

鉢头摩 $10^{4436777100798802905238461218816}$

僧祇 $10^{8873554201597605810476922437632}$

趣 $10^{17747108403195211620953844875264}$

至 $10^{35494216806390423241907689750528}$

阿僧祇 $10^{70988433612780846483815379501056}$

阿僧祇转 $10^{141976867225561692967630759002112}$

无量 $10^{283953734451123385935261518004224}$

无量转 $10^{567907468902246771870523036008448}$

无边 $10^{1135814937804493543741046072016896}$

无边转 $10^{2271629875608987087482092144033792}$

无等 $10^{4543259751217974174964184288067584}$

无等转 $10^{9086519502435948349928368576135168}$

不可数 $10^{18173039004871896699856737152270336}$

不可数转 $10^{36346078009743793399713474304540672}$

不可称 $10^{72692156019487586799426948609081344}$

不可称转 $10^{145384312038975173598853897218162688}$

不可思 $10^{290768624077950347197707794436325376}$

不可思转 $10^{581537248155900694395415588872650752}$

不可量 $10^{1163074496311801388790831177745301504}$

不可量转 $10^{2326148992623602777581662355490603008}$

不可说 $10^{4652297985247205555163324710981206016}$

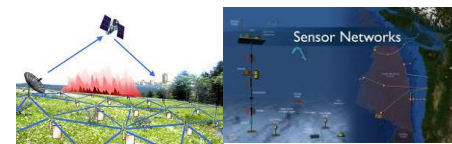
不可说转 $10^{9304595970494411110326649421962412032}$

不可说不可说 $10^{1860919194098882220653298843924824064}$

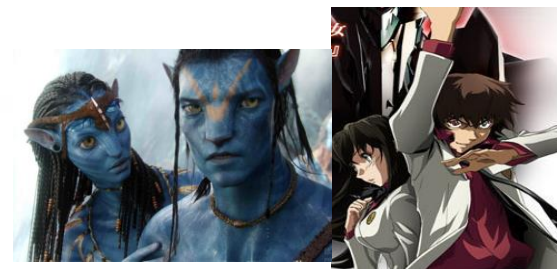
不可说不可说转 $10^{37218383881977644441306597687849648128}$

不可说不可说不可说 $10^{74436767763955288882613195375699296256}$

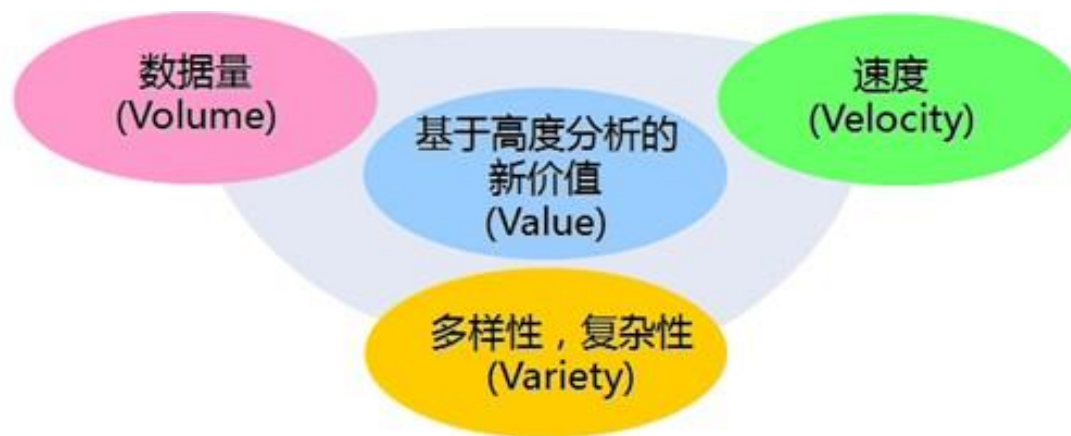
处处皆是大数据



商业数据



大数据的特点



大数据的应用

- 预测
- 推荐
- 商业情报分析
- 科学研究

应用中的大数据算法

- 预测
- 推荐
- 商业情报分析
- 科学研究

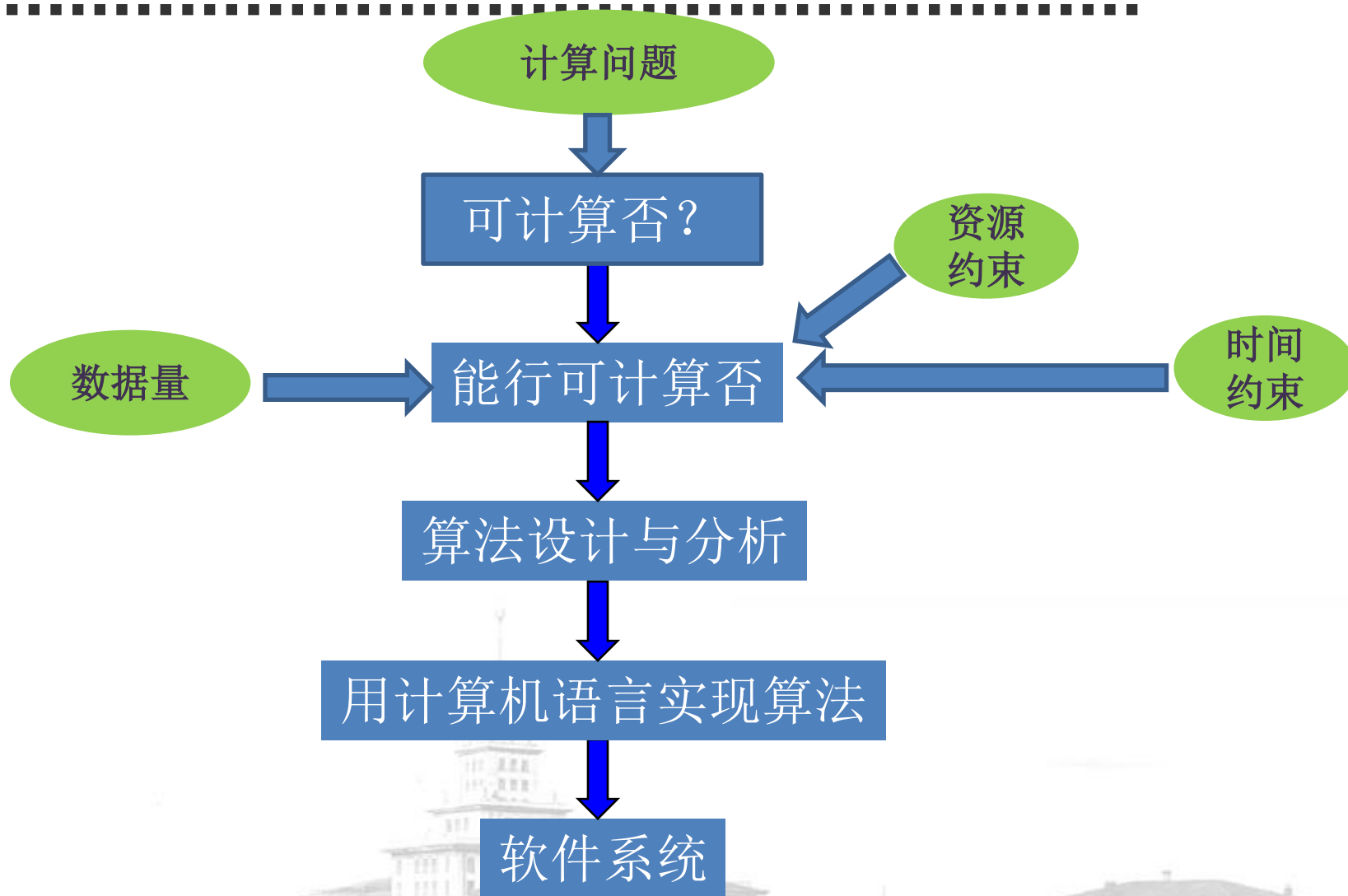
本讲内容

1.1 大数据的定义与特点

1.2 大数据算法

1.3 大数据算法设计与分析

大数据上问题求解计算问题的过程



大数据算法

● 大数据算法的定义

- 在给定的**资源约束**下，以大数据为输入，在给定的**时间约束**内可以生成满足给定**约束**结果的算法。

● 大数据算法可以不是：

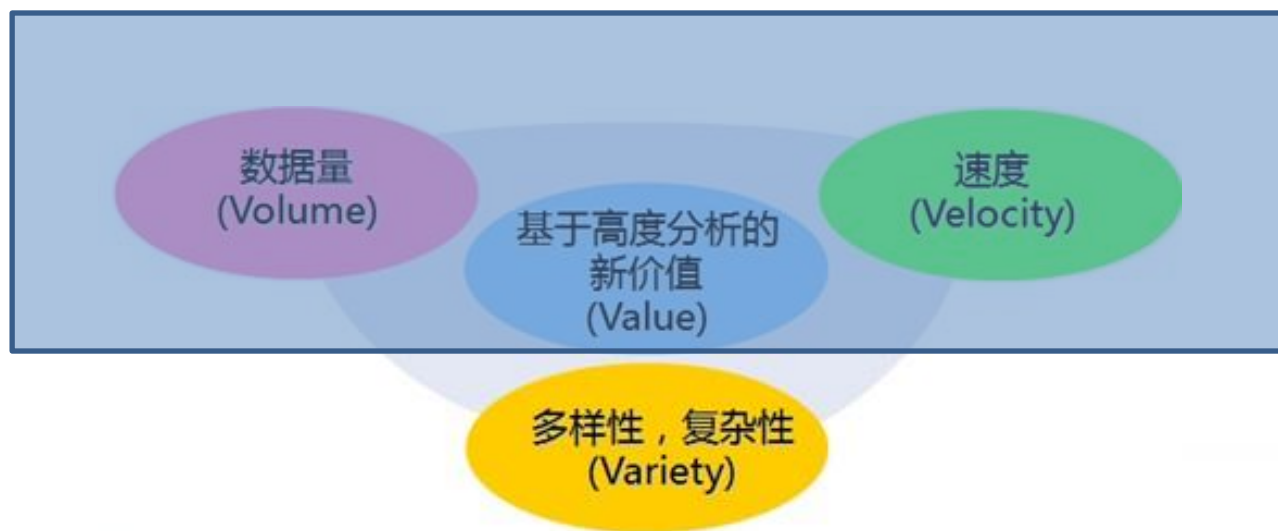
- 精确算法
- 内存算法
- 串行算法
- 仅在电子计算机上运行的算法

和“算法设计与分析”
课程中算法大不相同

● 大数据算法不仅是：

- 云计算
- MapReduce
- 大数据分析和挖掘的算法
- 数据库中的算法

大数据的特点与大数据算法



大数据算法的难度

- 访问全部数据时间过长

- 读取部分数据

时间亚线性算法

- 数据难于放入内存计算

- 将数据存储到磁盘上
- 仅基于少量数据进行计算

外存算法

空间亚线性算法

- 单个计算机难以保存全部数据，计算需要整体数据

- 并行处理

并行算法

- 计算机计算能力不足或知识不足

- 人来帮忙

众包算法

本讲内容

1.1 大数据的定义与特点

1.2 大数据算法

1.3 大数据算法设计与分析

大数据的算法设计技术

- 精确算法设计方法
- 并行算法
- 近似算法
- 随机算法
- 在线算法/数据流算法
- 外存算法
- 面向新型体系结构的算法
- 现代优化算法



大数据的算法分析

- 时间空间复杂性
- IO复杂性
- 结果质量（近似比、competitive ratio）
- 通讯复杂性



这门课的内容

- 亚线性算法
- 外存算法
- 并行算法
- 众包算法

