



大数据处理系统简介

金熠波

2019年5月15日

逸夫楼C-115, 18:30



提纲

- 大数据处理系统**发展**
 - 集群化处理 → 基于内存的数据处理
- 大数据处理系统**应用**
 - MapReduce型范式 → DAG型范式
- 大数据处理系统**剖析**
 - Hadoop → Spark
- **异构**硬件加速大数据系统



大数据处理系统发展

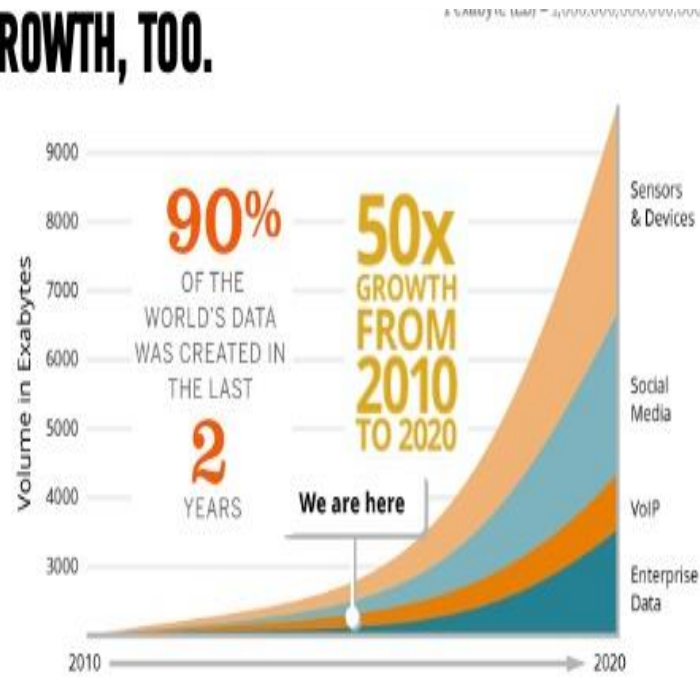


大数据处理需求兴起

- Google: 大规模网页数据
 - 需求目标: 存储和索引 (Page Rank)



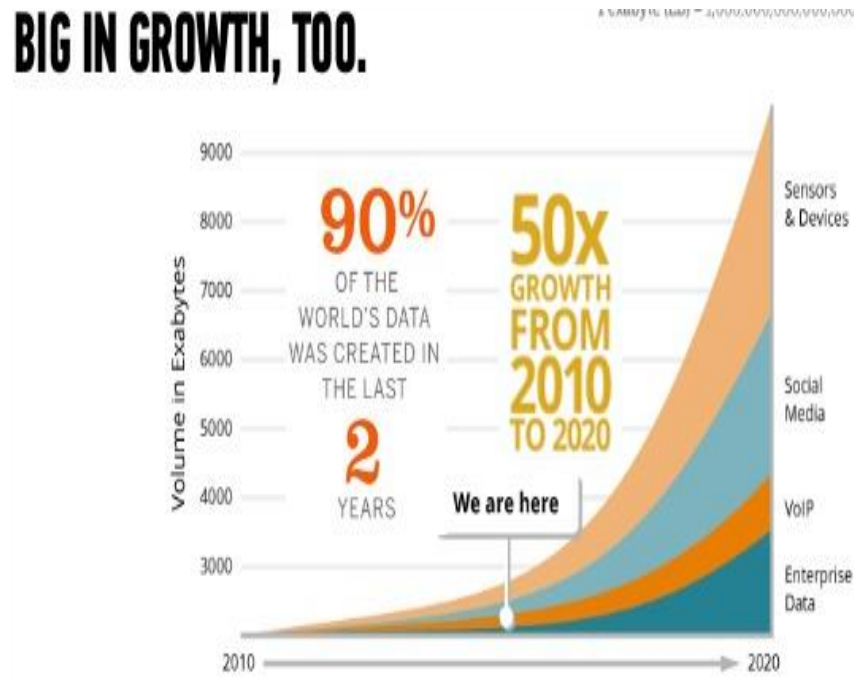
BIG IN GROWTH, TOO.





大数据处理需求兴起

- Google: 大规模网页数据
 - 需求目标: 存储和索引 (Page Rank)





集群化数据处理

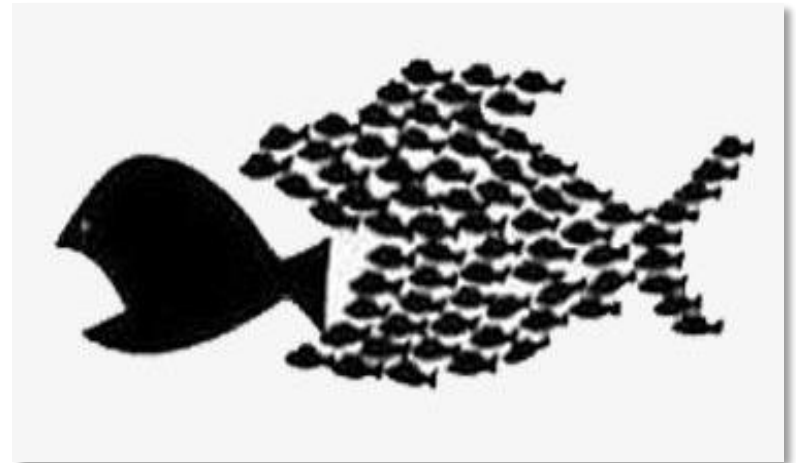
- Google: 大规模网页数据
 - 需求目标: 存储和索引
 - 解决方案: 利用**集群**并行化处理数据 OSDI' 04
 - Apache Hadoop '08: MapReduce



整合集群资源



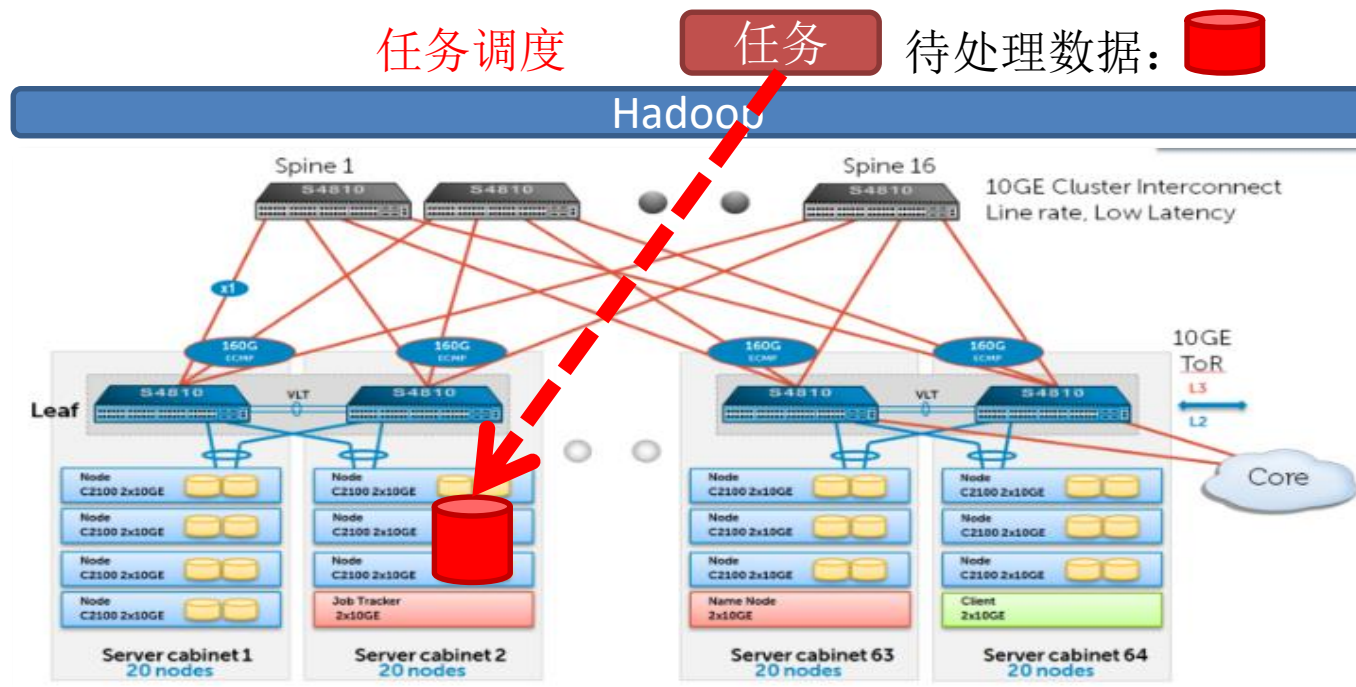
任务并发执行





传统集群内大数据处理

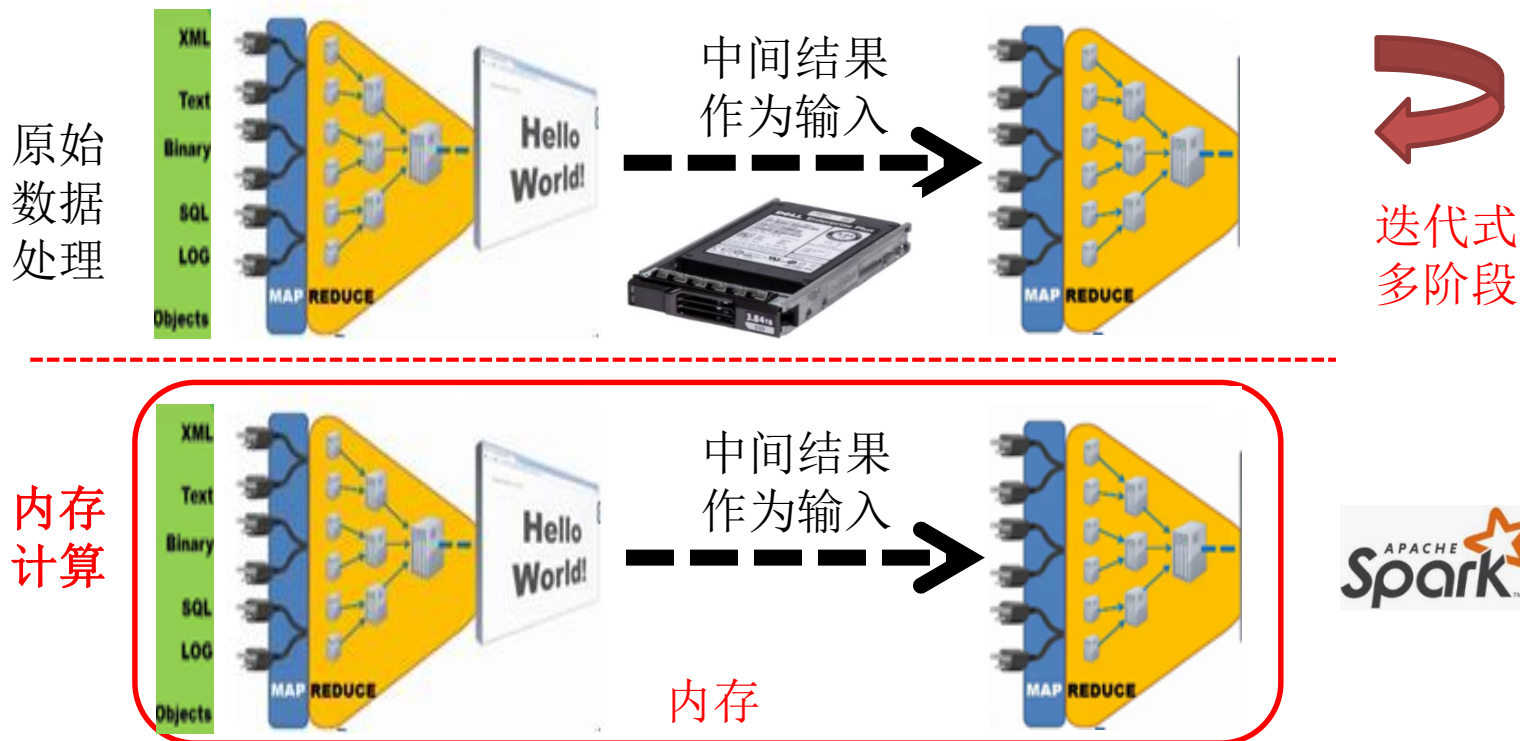
- 核心思想：利用数据**本地性**进行处理
 - 避免进行大规模数据移动Eurosys' 10





集群大数据处理发展

- Spark: 避免频繁的磁盘操作
 - 基于内存的分布式弹性数据集RDD, NSDI' 12





跨域大数据处理需求

- 全球化业务的发展促进跨域数据处理需求
 - 基于区域的广告词竞价: Google Adwords
 - 跨域视频流处理: 阿里云、优酷视频点播

Google Search-based Keyword Tool

Keywords | My draft keywords (0)

Website: nanjingb2c.com | Word or phrase: B2C | Advanced Search | Find keywords | Or browse all keywords

Categories for this search: All categories

Help: Learning the basics, Filtering and saving keywords, Incorporating into AdWords, Troubleshooting

Keyword ideas for (China) edit

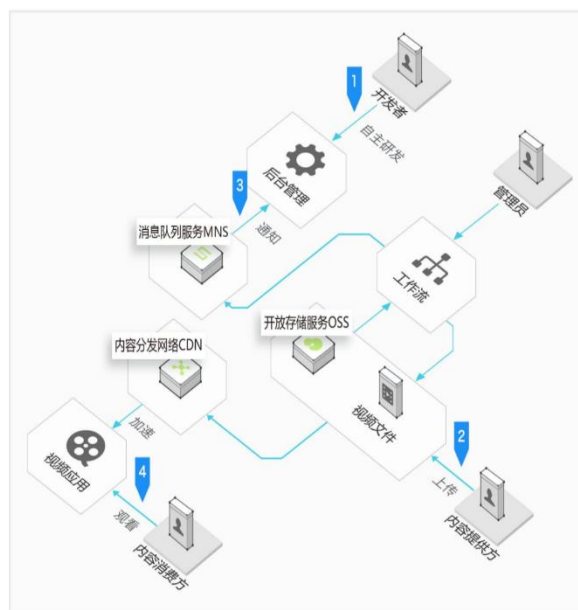
Our system doesn't have enough recent history on your account to provide meaningful suggestions.

Save to draft | Export

1-17 of 17

Keyword	Monthly searches	Competition	Sugg. bid	Ad Search share	Extracted from website
Keywords related to nanjingb2c.com (0)					
No keyword ideas for this site matching the search criteria and category.					
Keywords related to your search words and phrases (17)					
b2c电子商务解决方案	125		CNY2.17	-	-
b2c电子商务平台	100		CNY4.65	-	-
b2c平台	82		CNY2.48	-	-
b2c购物	52		CNY2.43	-	-
外贸b2c	42		CNY4.68	-	-
百货b2c	42		CNY8.63	-	-
中国b2c	42		CNY10.22	-	-

51CTO.COM 技术成就梦想



视频点播

更流畅点播体验

提供集音视频上传、自动化转码、媒体资源管理、分发加速于一体的一站式音视频点播解决方案。帮助快速搭建安全、弹性、高可定制的点播平台和应用程序。

能够解决

灵活自动专业的媒体处理

零开发自定义云端音视频处理流程，可视化配置转码、截图、水印...

访问控制，版权保护

提供referer黑白名单和URL鉴权，彻底解决盗链危害

播放流畅无卡顿

稳定高效的性能：95%+命中率，毫秒级响应时间

推荐搭配使用

媒体转码 OSS 消息服务



跨域集群现状

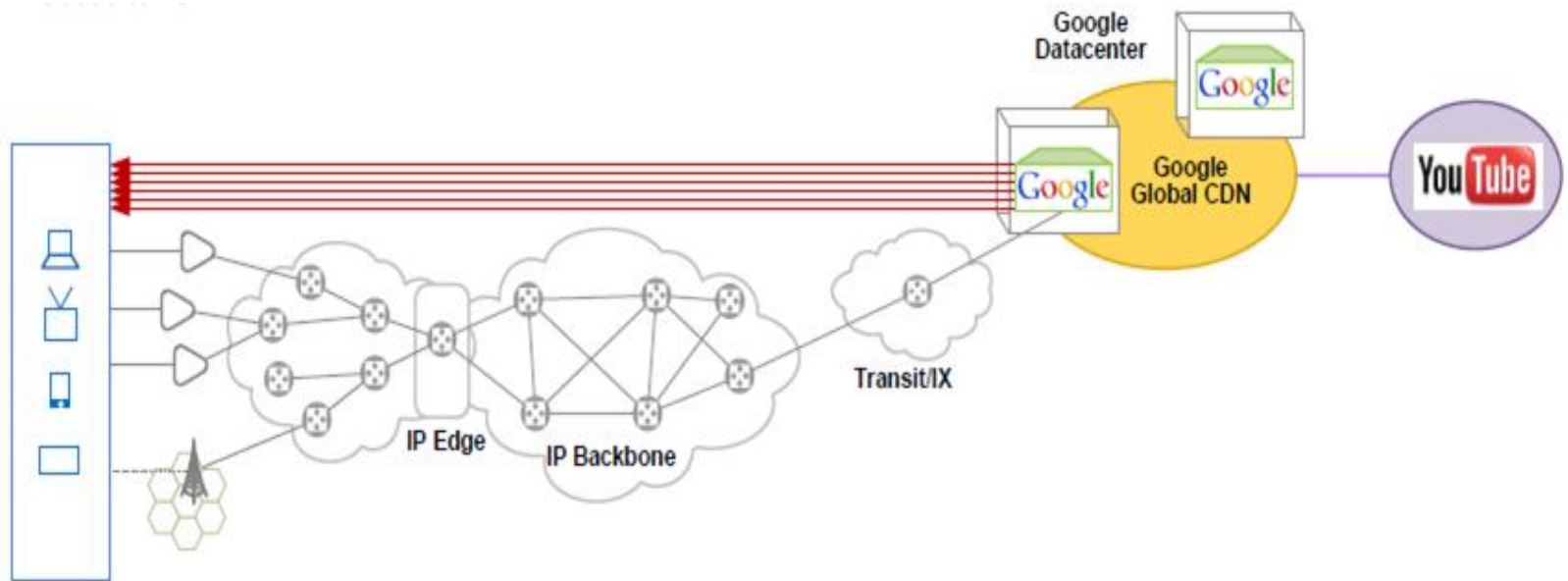
- 数据中心向边缘集群的演变





跨域集群现状

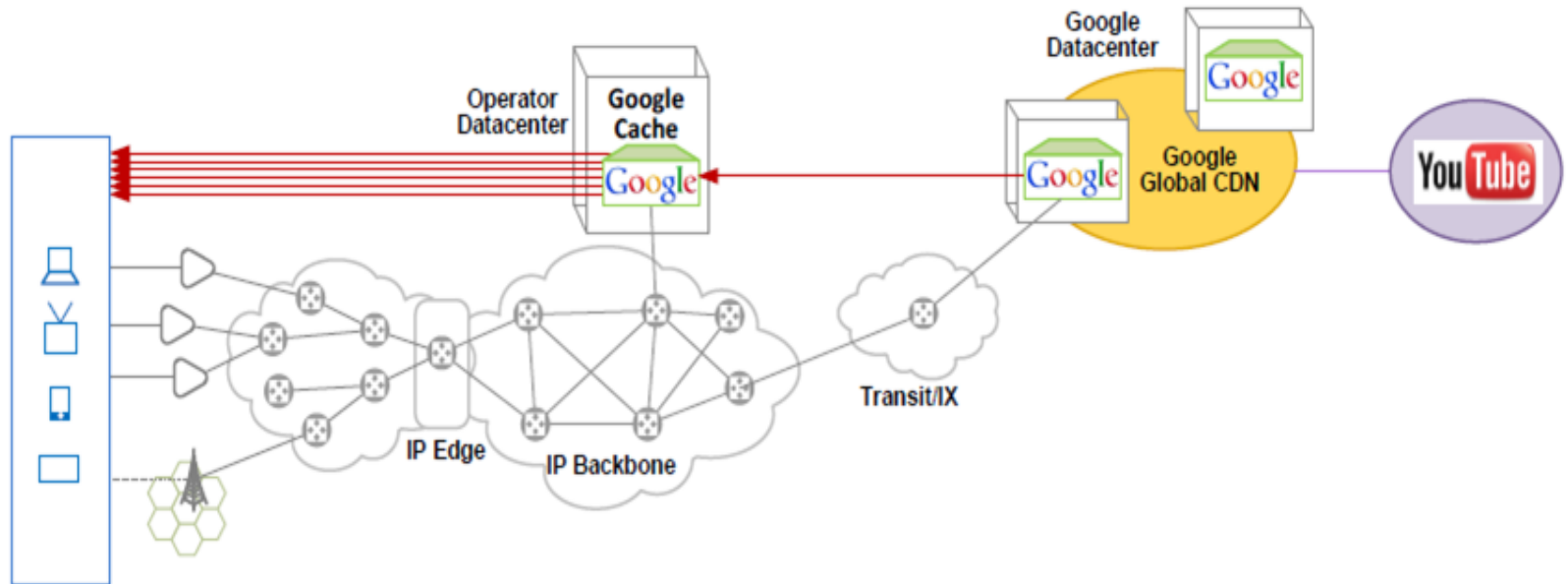
- 数据中心向边缘集群的演变
 - 跨域响应时间长



<https://www.netmanias.com/en/post/blog/5921/cdn-google-iptv-netflix-ott-video-streaming-youtube/who-wins-ott-cache-vs-operator-cache>

跨域集群现状

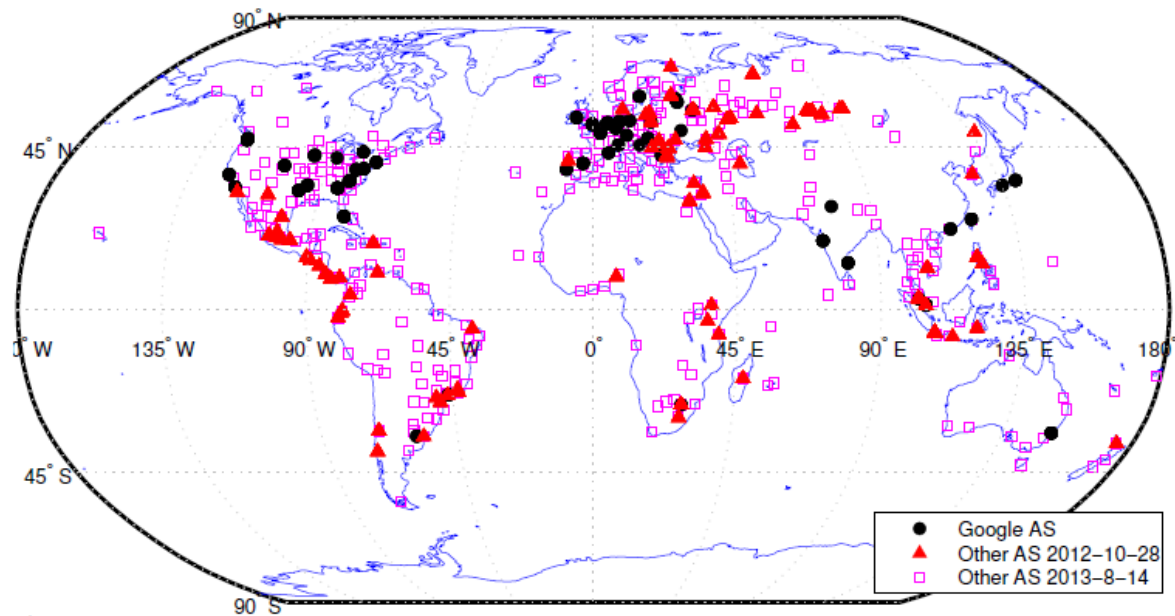
- 数据中心向边缘集群的演变
 - 互联网服务提供商部署边缘处理节点



<https://www.netmanias.com/en/post/blog/5921/cdn-google-iptv-netflix-ott-video-streaming-youtube/who-wins-ott-cache-vs-operator-cache>

跨域集群现状

- 数据中心向边缘集群的演变
 - 数以千计的边缘集群协同数据中心提供服务





Akamai Platform Architecture

Jan. 2013

- Mapping System (centralized)
- HL-DNS (centralized)
- Monitoring Agent (distributed/located in Cluster)
- LL-DNS (distributed/located in Cluster)
- Edge Server (distributed, totally 119,000 servers)

The diagram illustrates the Akamai Platform Architecture. It shows a central Mapping System (Map Maker) that updates a High Level DNS Mapping Table (HL-DNS) and Low Level DNS Mapping Tables (LL-DNS) across various clusters (Cluster 1, Cluster 2, ..., Cluster 1,000, ..., Cluster 5,000+). Each cluster contains Monitoring Agents and Edge Servers. The Monitoring Agents report periodically to the Mapping System and check server load and health. The Mapping System updates the HL-DNS and LL-DNS tables. The LL-DNS tables are used by the Edge Servers to deliver content to the user. The user's local DNS server is also shown, which receives updates from the Mapping System and the Edge Servers. The diagram includes a legend for the components and a list of three key update processes: 1. Reporting Periodically (Monitoring Agent to Mapping System), 2. Low Level Map Update (Mapping System to Low-Level DNS), and 3. High Level Map Update (Mapping System to High-Level DNS).

1 Reporting Periodically (Monitoring Agent to Mapping System)

- Health & Load of Clusters and Edge Servers (a)
- RTT & Packet Loss between Clusters and Local DNS Servers (b)
- RTT & Packet Loss between Akamai Clusters

2 Low Level Map Update: Every 2~10s (Mapping System to Low-Level DNS)

- Edge Server Status in a Cluster: Health & Load of Edge Servers
- RTT & Packet Loss between Clusters and Local DNS Servers

3 High Level Map Update: Every 15~20m (Mapping System to High-Level DNS)

- Mapping between LL-DNS Servers and Local DNS Servers
- Cluster Status: Health & Load of Cluster
- RTT & Packet Loss between Clusters and Local DNS Servers

www.netmanias.com

Netmanias ONE-SHOW



边缘集群产生大量数据

- 例如：用于商业广告竞价的点击数据
 - Adwords：利用该数据进行广告竞价和收费

```
ClickLog(sourceIP, destURL, visitDate, adRevenue, ...)
```

```
Q: SELECT sourceIP, sum(adRevenue), avg(pageRank)
   FROM ClickLog cl JOIN PageInfo pi
      ON cl.destURL = pi.pageURL
  WHERE pi.pageCategory = 'Entertainment'
  GROUP BY sourceIP
  HAVING sum(adRevenue) >= 100
```



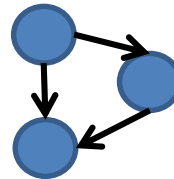


跨域大数据处理框架

- 跨域数据处理架构

```
Q: SELECT sourceIP, sum(adRevenue), avg(pageRank)
FROM ClickLog cl JOIN PageInfo pi
  ON cl.destURL = pi.pageURL
WHERE pi.pageCategory = 'Entertainment'
GROUP BY sourceIP
HAVING sum(adRevenue) >= 100
```

数据分析查询

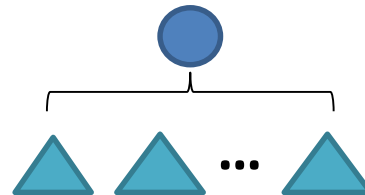


数据分析作业

Frameworks: Spark/Hadoop



跨域数据集



批次数据分析任务

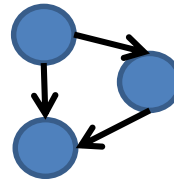


跨域大数据处理框架

- 跨域数据处理架构

```
Q: SELECT sourceIP, sum(adRevenue), avg(pageRank)
FROM ClickLog cl JOIN PageInfo pi
  ON cl.destURL = pi.pageURL
WHERE pi.pageCategory = 'Entertainment'
GROUP BY sourceIP
HAVING sum(adRevenue) >= 100
```

数据分析查询

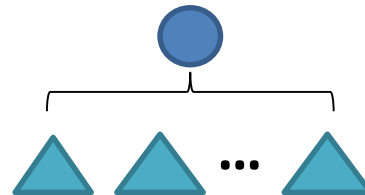


数据分析作业

Frameworks: Spark/Hadoop



跨域数据集

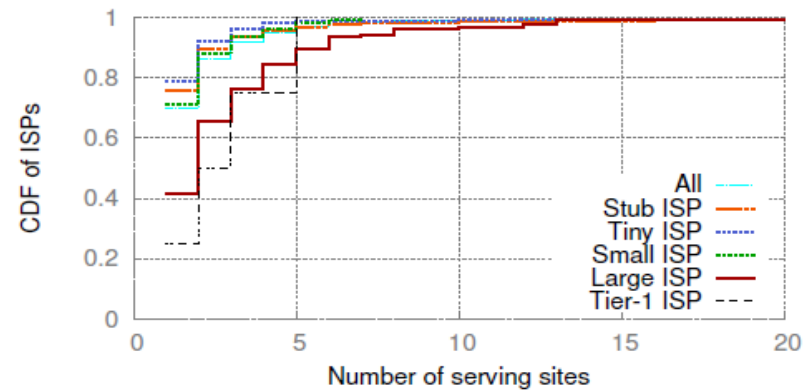
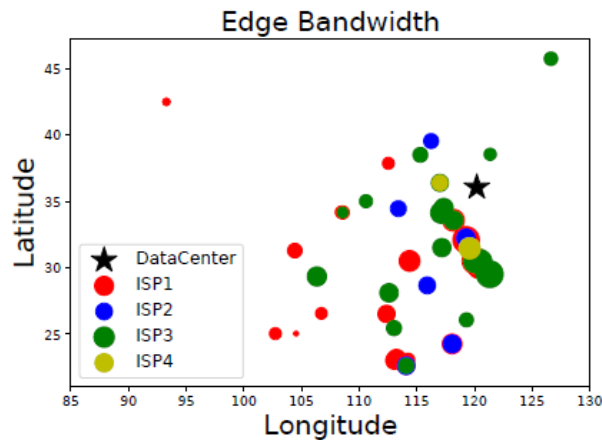


减少批次任务
完成时间



跨域边缘集群的异构性

- 资源异构性
 - 带宽资源
 - 计算资源



Performance		Worst	Best
Server	Memory	GB	TB, DDR
	CPU (Cores)	Few	Hundreds
	Disk	GB	TB, SSD

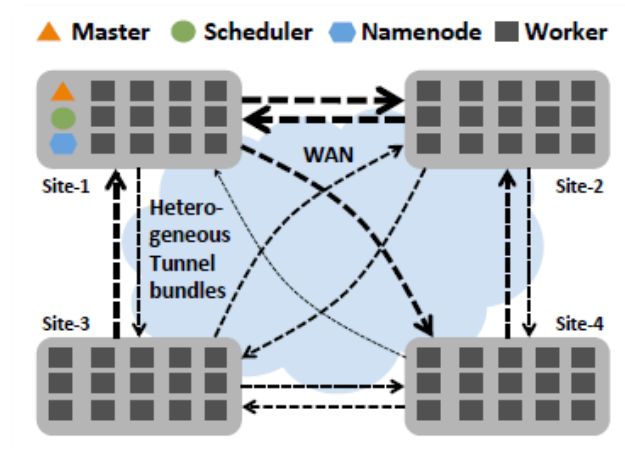
Calder *et. al.* Mapping the Expansion of Google's Serving Infrastructure. IMC' 13
Alibaba Cloud Computing, and Chinaz for measurement.



集群间数据处理

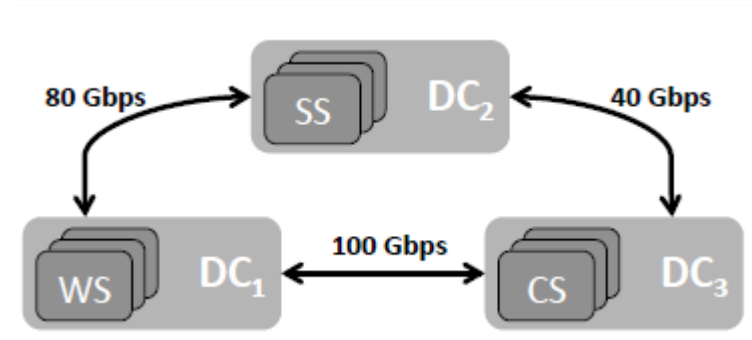
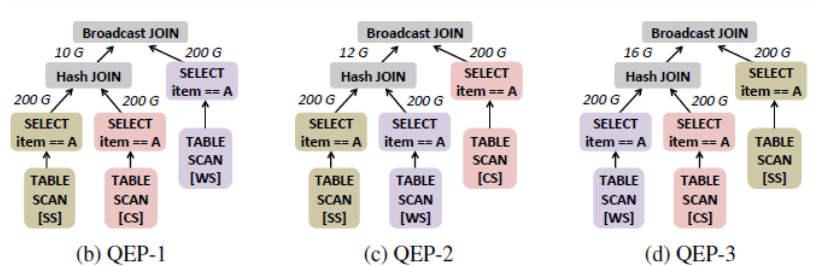
内存计算

- 跨域资源异构性大
 - 需要综合计算力 & 带宽进行资源部署



Performance		Worst	Best
Server	Memory	GB	TB, DDR3/4
	CPU	Several Cores	Hundreds Cores
	Disk	GB	TB, SSD

(b) Edges with heterogeneous computing capacities.





大数据处理生态



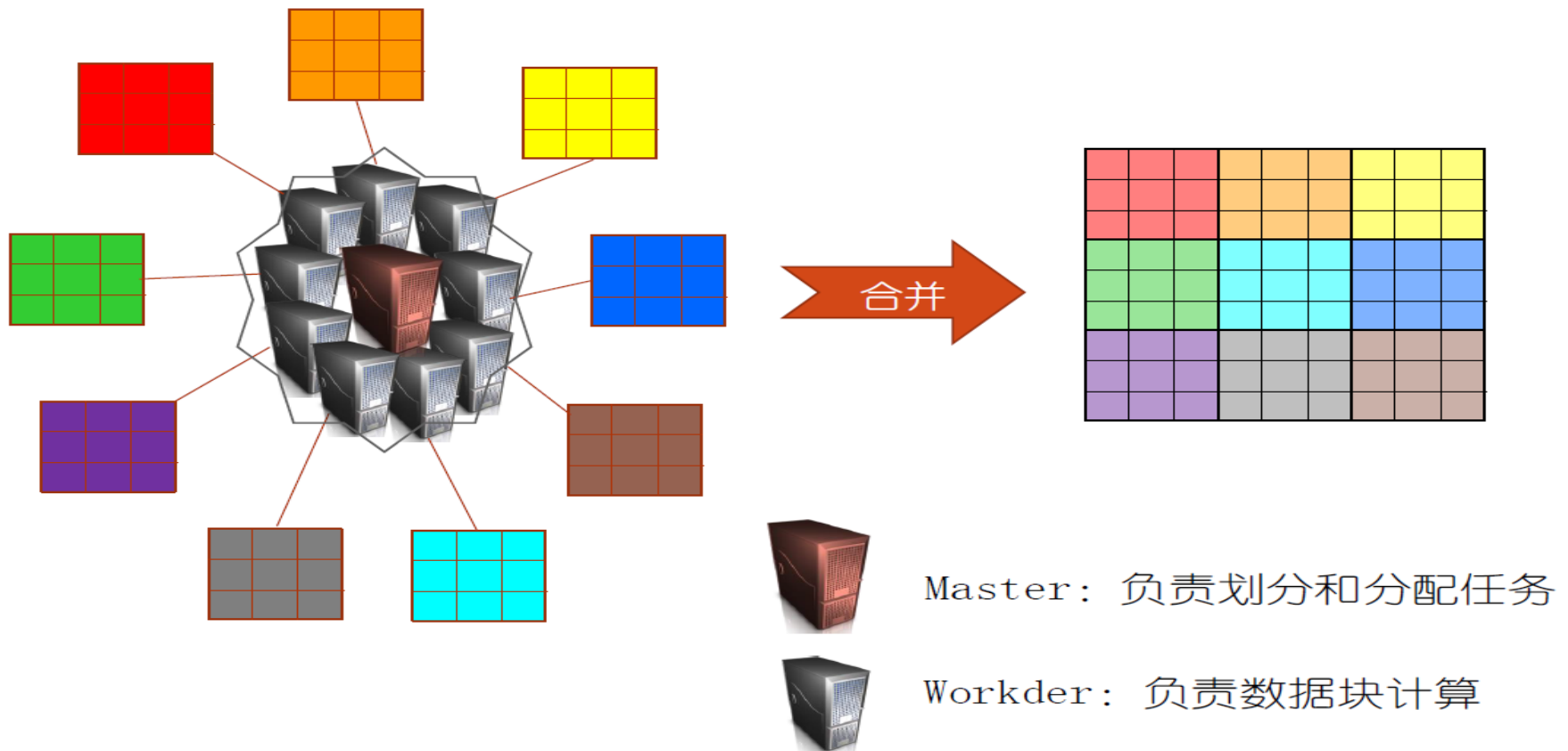


大数据处理系统应用



大数据处理应用

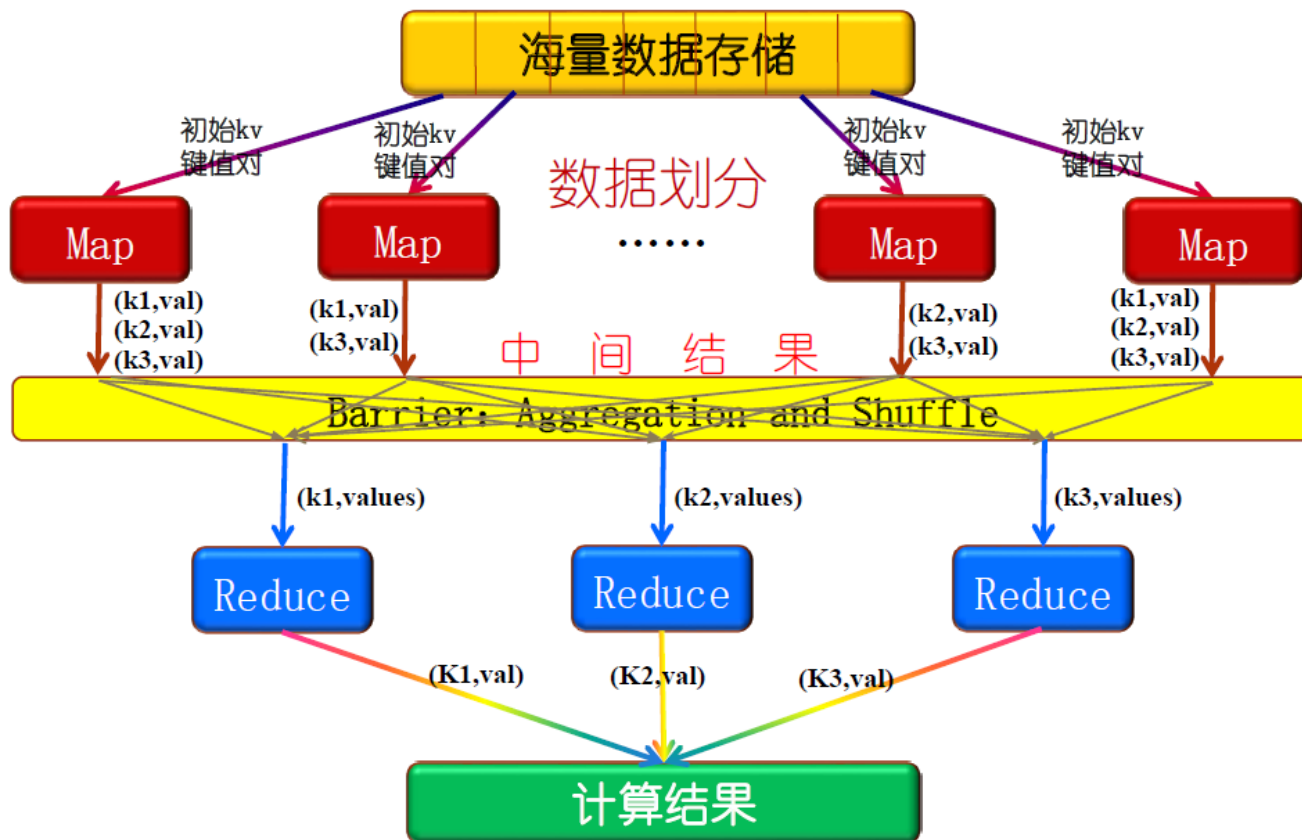
- MapReduce范式





大数据处理应用

- MapReduce范式





大数据处理应用

基于MapReduce的处理过程示例—文档词频统计：WordCount

设有4组原始文本数据：

Text 1: the weather is good

Text 2: today is good

Text 3: good weather is good

Text 4: today has good weather

传统的串行处理方式(Java):

```
String[] text = new String[]
{ "the weather is good", "today is good", "good weather is good", "today has good weather" };
HashTable ht = new HashTable();
for(i=0; i<3; ++i)
{ StringTokenizer st = new StringTokenizer(text[i]);
  while (st.hasMoreTokens())
  { String word = st.nextToken();
    if(!ht.containsKey(word)) { ht.put(word, new Integer(1)); }
    else { int wc = ((Integer)ht.get(word)).intValue() +1; // 计数加1
          ht.put(word, new Integer(wc));
        }
  }
}
for (Iterator itr=ht.keySet().iterator(); itr.hasNext(); )
{ String word = (String)itr.next();
  System.out.print(word+ ": " + (Integer)ht.get(word)+" ";
}
```

输出: good: 5; has: 1; is: 3; the: 1; today: 2; weather: 3



大数据处理应用

构建抽象模型：Map与Reduce

基于MapReduce的处理过程示例—文档词频统计：WordCount

MapReduce处理方式

使用4个map节点：

map节点1:

输入：(text1, "the weather is good")

输出：(the, 1), (weather, 1), (is, 1), (good, 1)

map节点2:

输入：(text2, "today is good")

输出：(today, 1), (is, 1), (good, 1)

map节点3:

输入：(text3, "good weather is good")

输出：(good, 1), (weather, 1), (is, 1), (good, 1)

map节点4:

输入：(text3, "today has good weather")

输出：(today, 1), (has, 1), (good, 1), (weather, 1)



大数据处理应用

构建抽象模型：Map与Reduce

基于MapReduce的处理过程示例—文档词频统计：WordCount

MapReduce处理方式

使用3个reduce节点：

reduce节点1:

输入：(good, 1), (good, 1), (good, 1), (good, 1), (good, 1)

输出：(good, 5)

reduce节点2:

输入：(has, 1), (is, 1), (is, 1), (is, 1),

输出：(has, 1), (is, 3)

reduce节点3:

输入：(the, 1), (today, 1), (today, 1)

(weather, 1), (weather, 1), (weather,

输出：(the, 1), (today, 2), (weather, 3)

输出：

good: 5

is: 3

has: 1

the: 1

today: 2

weather: 3



大数据处理应用

构建抽象模型：Map与Reduce

基于MapReduce的处理过程示例—文档词频统计：WordCount

MapReduce处理方式

MapReduce伪代码(实现Map和Reduce两个函数)：

```
Class WordCountMapper
  method map(String input_key, String input_value):
    // input_key: text document name
    // input_value: document contents
    for each word w in input_value:
      EmitIntermediate(w, "1");
```

```
Class WordCountReducer
  method reduce(String output_key,
                Iterator intermediate_values):
    // output_key: a word
    // output_values: a list of counts
    int result = 0;
    for each v in intermediate_values:
      result += ParseInt(v);
    Emit(output_key, result);
```



大数据处理应用

构建抽象模型：Map与Reduce

基于MapReduce的处理过程示例—文档词频统计：WordCount

MapReduce处理方式

使用3个reduce节点：

reduce节点1:

输入：(good, 1), (good, 1), (good, 1), (good, 1), (good, 1)

输出：(good, 5)

reduce节点2:

输入：(has, 1), (is, 1), (is, 1), (is, 1),

输出：(has, 1), (is, 3)

reduce节点3:

输入：(the, 1), (today, 1), (today, 1)

(weather, 1), (weather, 1), (weather,

输出：(the, 1), (today, 2), (weather, 3)

输出：

good: 5

is: 3

has: 1

the: 1

today: 2

weather: 3

has + the 共有几个？



大数据处理应用

- 自定义Reducer划分

定制Partitioner

程序员可以根据需要定制Partitioner来改变Map中间结果到Reduce节点的分区方式，并在Job中设置新的Partitioner

```
Class NewPartitioner extends HashPartitioner<K,V>
{ // override the method
    getPartition(K key, V value, int numReduceTasks)
    { term = key. toString().split(",")[0]; //<term, docid>=>term
      super.getPartition(term, value, numReduceTasks);
    }
}
```

并在Job中设置新的Partitioner：

```
Job. setPartitionerClass(NewPartitioner)
```



大数据处理应用

- 自定义键值对：key-value

doc1:
one fish
two fish

doc2:
red fish
blue fish

doc3:
one red bird



倒排索引:
one: doc1, doc3
fish: doc1, doc2
two: doc1
red: doc2, doc3
blue: doc2
bird: doc3

基于以上索引的搜索结果:

fish → doc1, doc2

red → doc2, doc3

red fish → doc2



大数据处理应用

- 自定义键值对：key-value

简单的文档倒排算法

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class InvertedIndexMapper extends Mapper<Text, Text, Text, Text>
{
    @Override
    protected void map(Text key, Text value, Context context)
        throws IOException, InterruptedException
    // default RecordReader: LineRecordReader; key: line offset; value: line string
    {
        FileSplit fileSplit = (FileSplit)context.getInputSplit();
        String fileName = fileSplit.getPath().getName();
        Text word = new Text();
        Text fileName_lineOffset = new Text(fileName+"#" +key.toString());
        StringTokenizer itr = new StringTokenizer(value.toString());
        for(; itr.hasMoreTokens(); )
        {
            word.set(itr.nextToken());
            context.write(word, fileName_lineOffset);
        }
    }
}
```

改进:map输出的key除了文件名, 还给出了该词所在行的偏移值:
格式: filename#offset

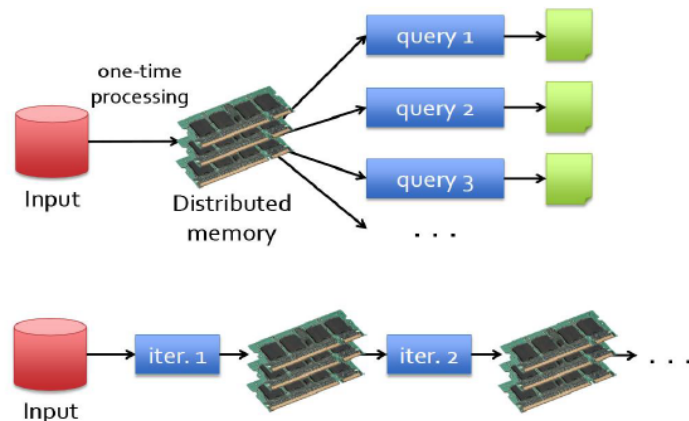
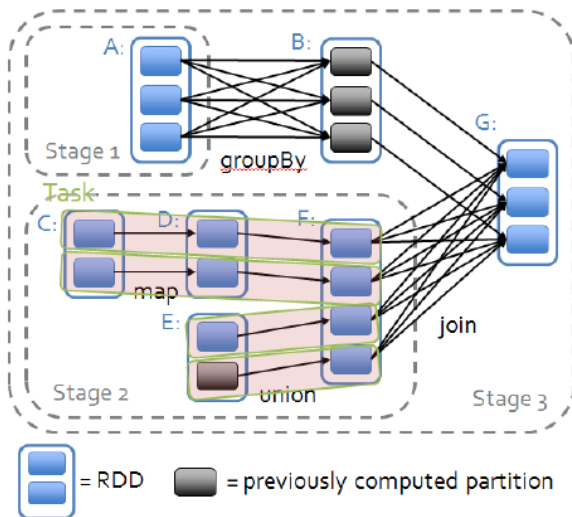


大数据处理应用

• DAG型计算范式：MapReduce的拓展

Spark基于内存计算思想提高计算性能

- Spark提出了一种基于内存的弹性分布式数据集(RDD)，通过对RDD的一系列操作完成计算任务，可以大大提高性能
- 同时一组RDD形成可执行的有向无环图DAG，构成灵活的计算流图
- 覆盖多种计算模式





大数据处理应用

- DAG型计算范式： MapReduce的拓展

Spark的基本编程方法与示例

//在一个存储于HDFS的Log文件中，计算出现ERROR的行数，本程序使用Scala语言编写，这个语言也是Spark开发和编程的推荐语言。

```
def main(args: Array[String]) { //定义一个main函数

    val conf = new SparkConf().setAppName("Spark Pi") //定义一个sparkConf，
        提供Spark运行的各种参数，如程序名称、用户名称等

    val sc = new SparkContext(conf) //创建Spark的运行环境，并将Spark运行的
        参数传入Spark的运行环境中

    val fileRDD=sc.textFile("hdfs:///root/Log") //调用Spark的读文件函数，从HDFS
        中读取Log文件，输出一个RDD类型的实例：fileRDD。具体类型：
        RDD[String]

    val filterRDD=fileRDD.filter(line=>line.contains("ERROR")) //调用RDD的filter函
        数，过滤fileRDD中的每一行，如果该行中含有ERROR，保留；否则，删
        除。生成另一个RDD类型的实例：filterRDD。具体类型:RDD[String]

    注： line=>line.contains("ERROR")表示对每一个line应用contains()函数

    val result = filterRDD.count() //统计filterRDD中总共有多少行，result为Int类型

    sc.stop() //关闭Spark
```



大数据处理应用

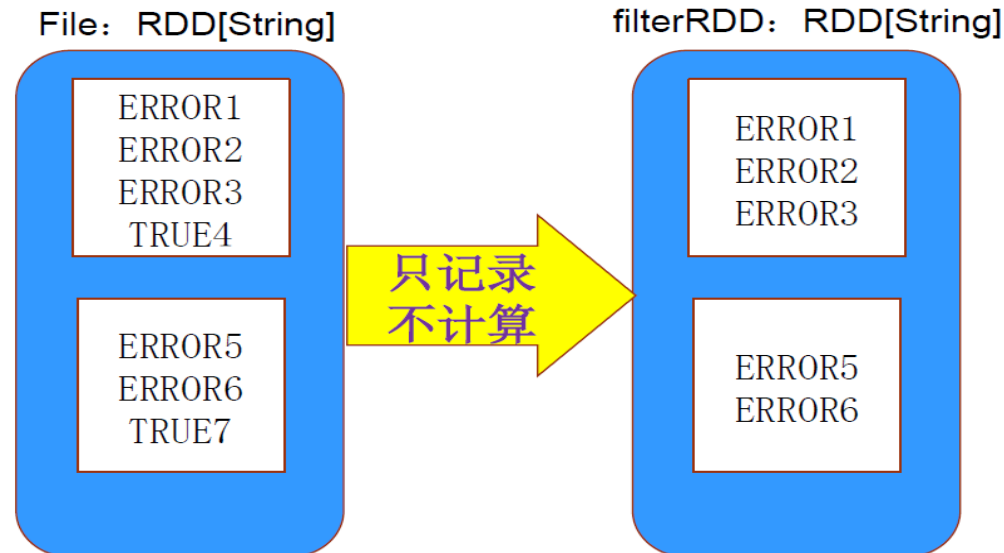
- DAG型计算范式：MapReduce的拓展

RDD的transformation示例

例： `val filterRDD=fileRDD.filter(line=>line.contains("ERROR"))`

设fileRDD中包含以下7行数据：

"ERROR1 ERROR2 ERROR3 TRUE4 ERROR5 ERROR6 TRUE7"



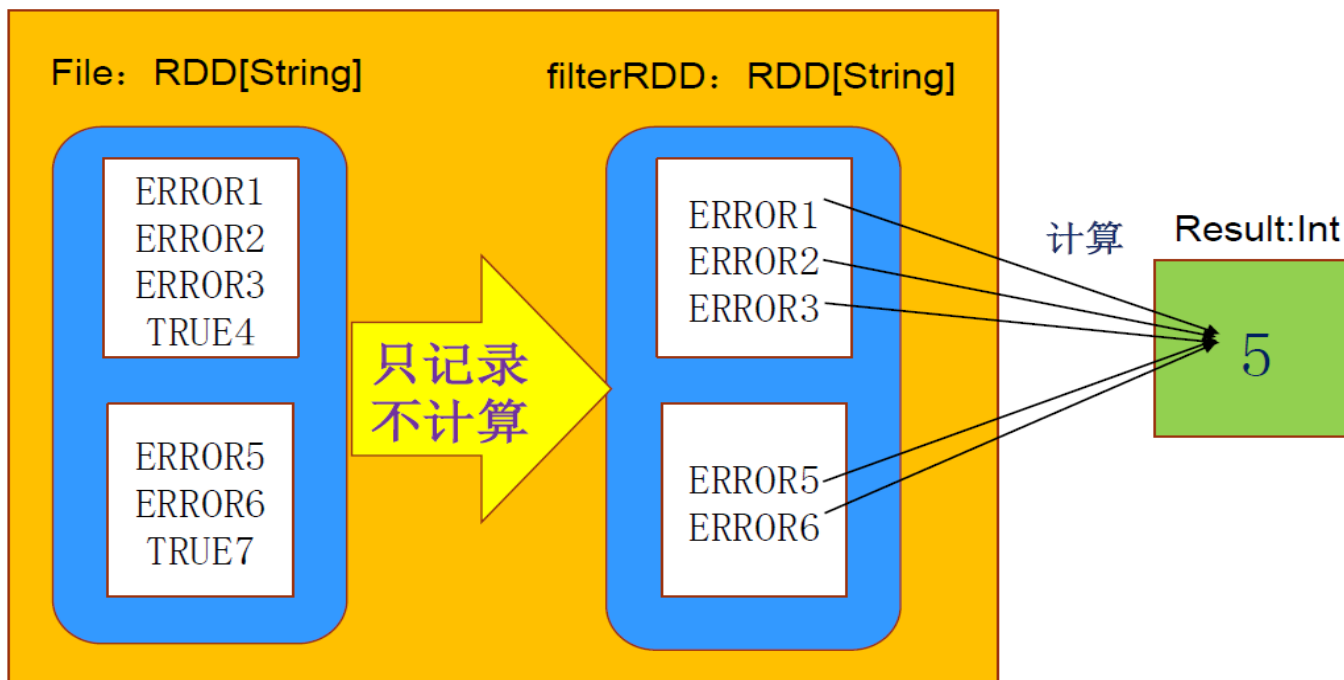


大数据处理应用

- DAG型计算范式：MapReduce的拓展

RDD的action图示

例： `val result = filterRDD.count()`





大数据处理应用

• DAG型计算范式：MapReduce的拓展

RDD之间的依赖关系

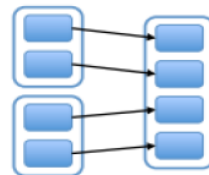
- 在Spark中存在两种类型的依赖：

- **窄依赖**：父RDD中的一个Partition最多被子RDD中的一个Partition所依赖。
- **宽依赖**：父RDD中的一个Partition被子RDD中的多个Partition所依赖。

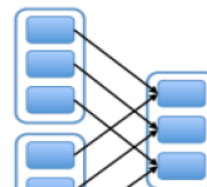
Narrow Dependencies:



map, filter



union

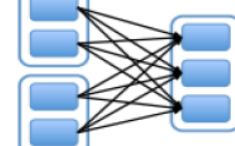


join with inputs
co-partitioned

Wide Dependencies:



groupByKey

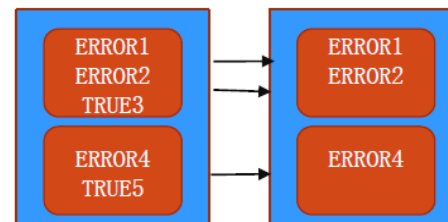


join with inputs not
co-partitioned

例子：

```
val filterRDD=fileRDD.filter  
                (line=>line.contains("ERROR"))
```

是一种窄依赖



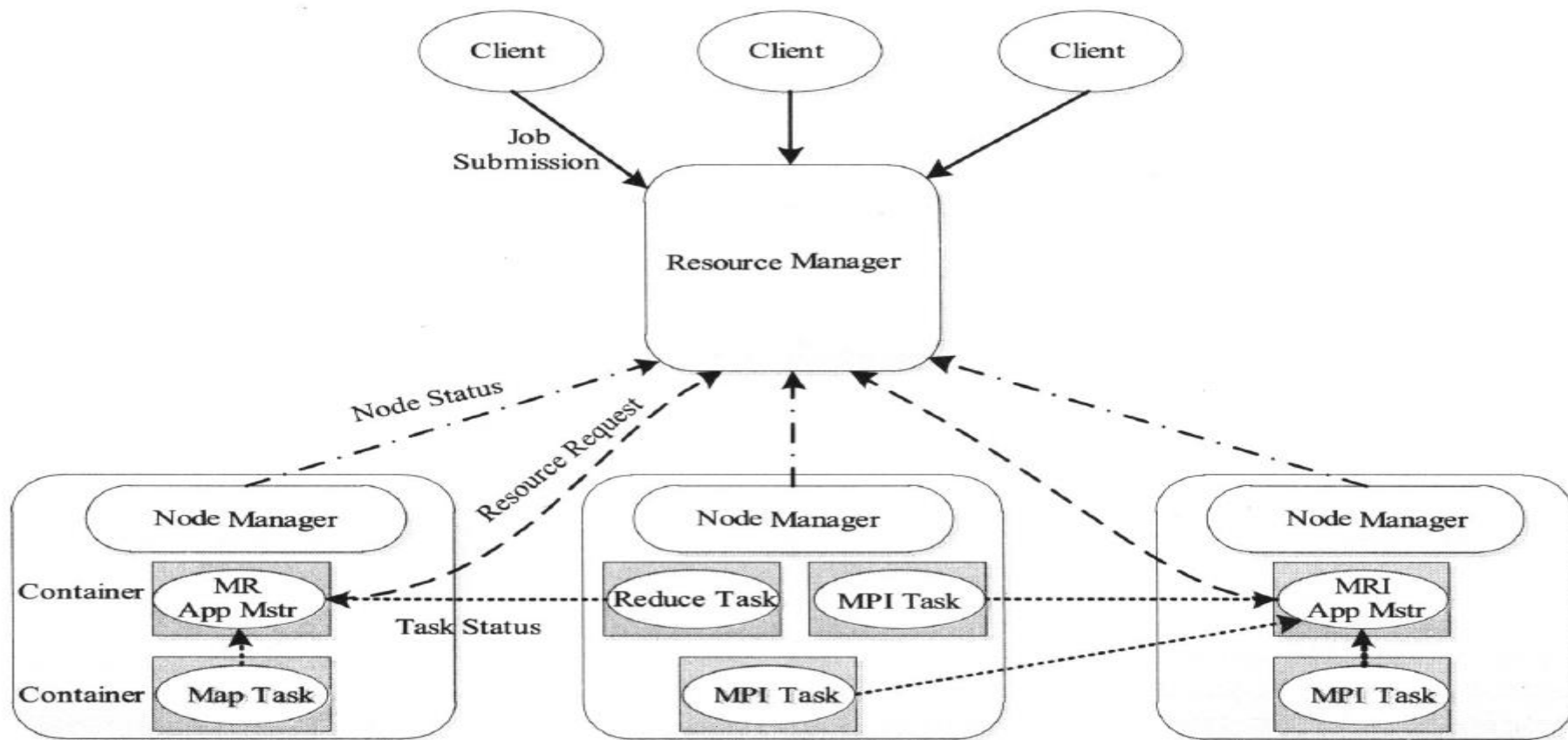


大数据处理系统剖析



大数据处理系统剖析

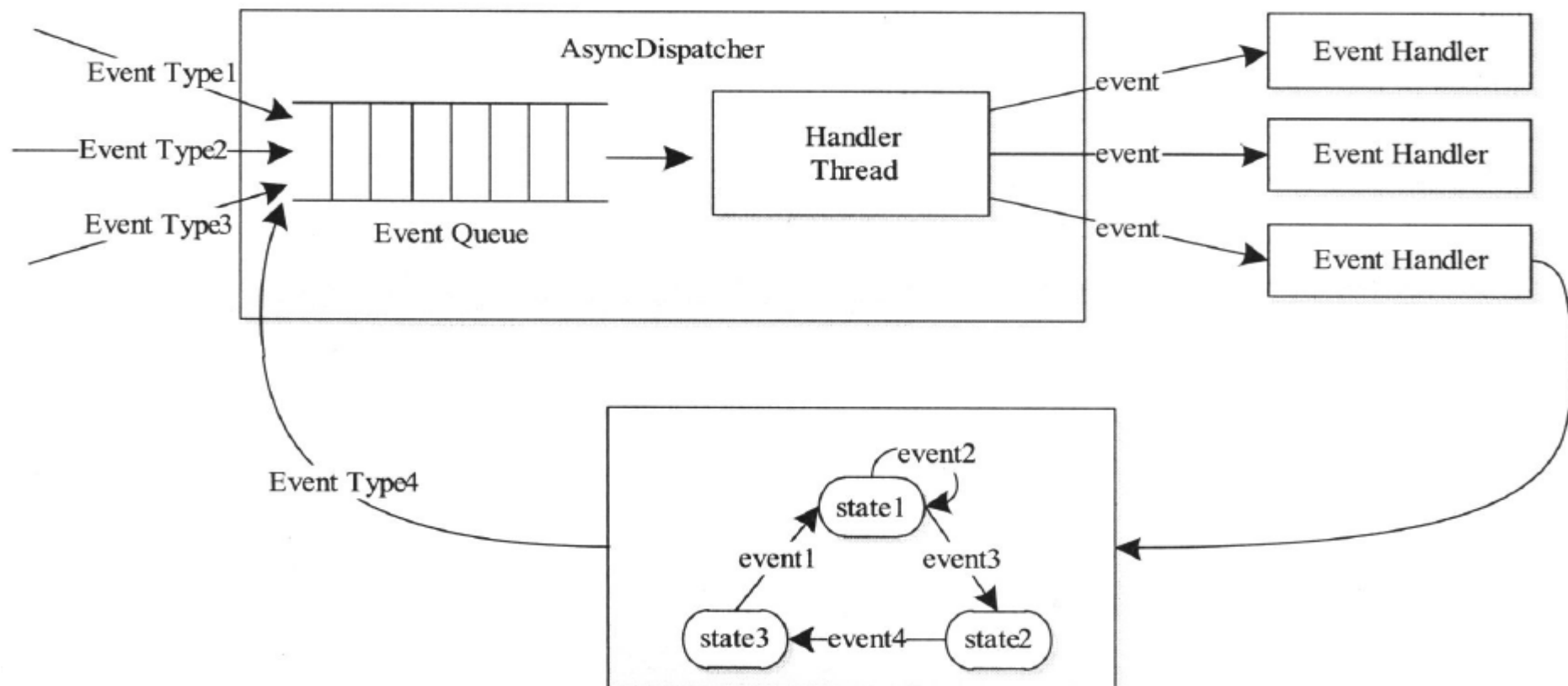
- 集群化资源管理Hadoop-Yarn





大数据处理系统剖析

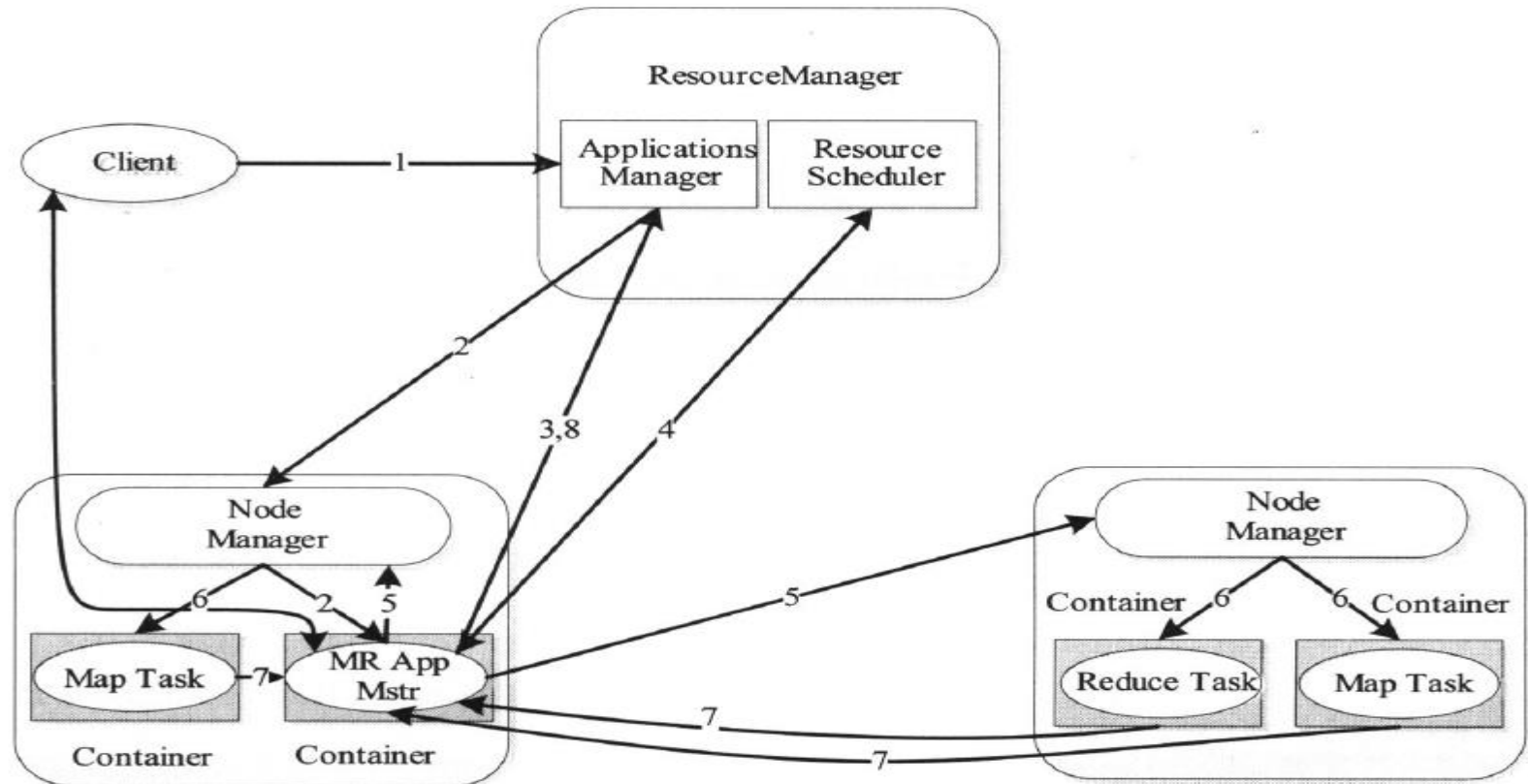
- Hadoop事件处理模块





大数据处理系统剖析

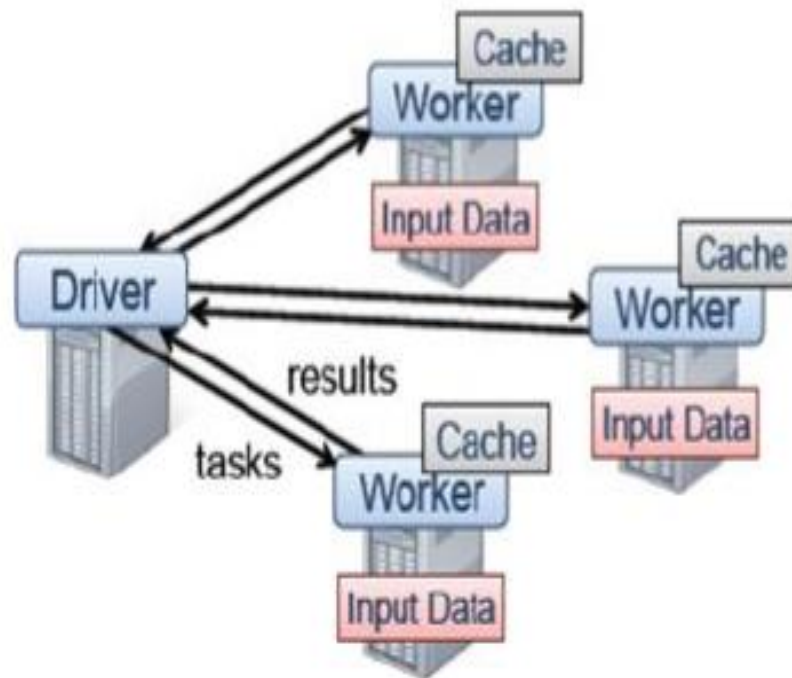
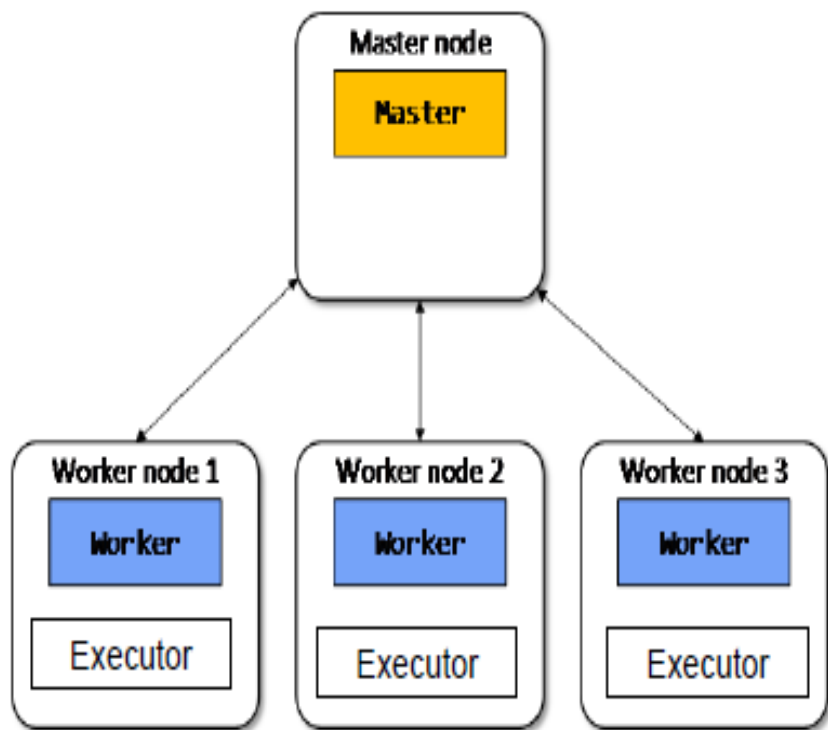
- Hadoop资源分配流程





大数据处理系统剖析

- 集群化资源管理Spark





异构硬件加速大数据系统



异构硬件数据处理背景

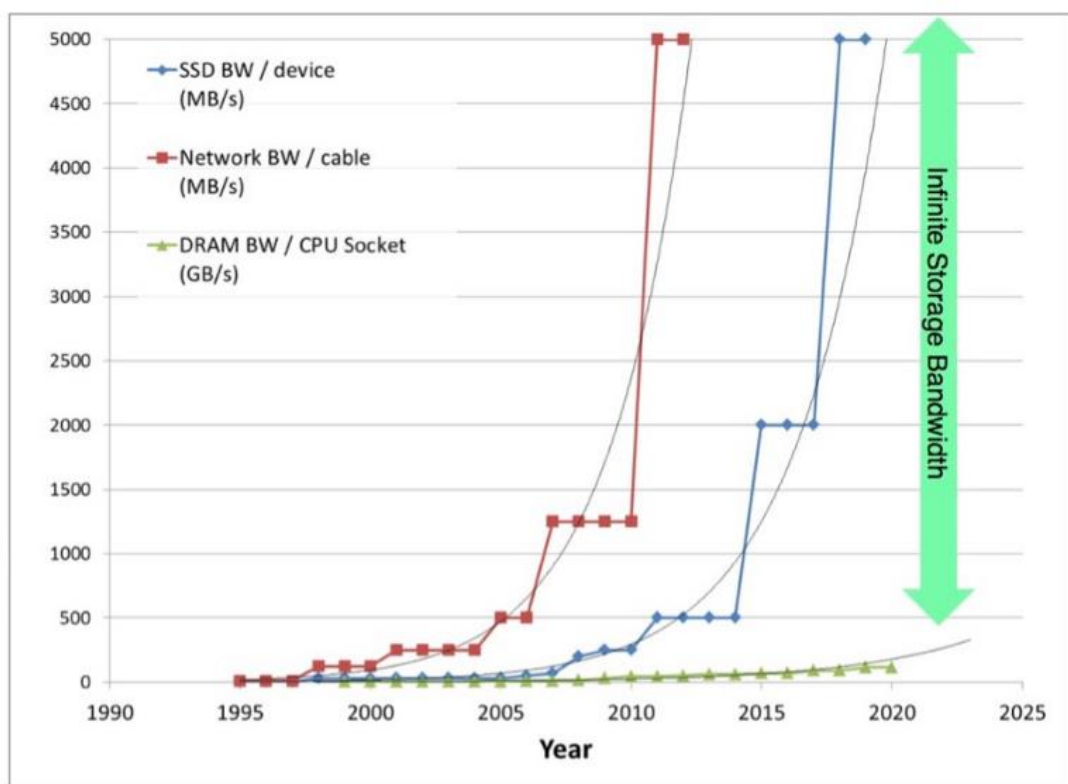
异构硬件

- 异构硬件突飞猛进

Network, Storage, & DRAM trends

Linear scale

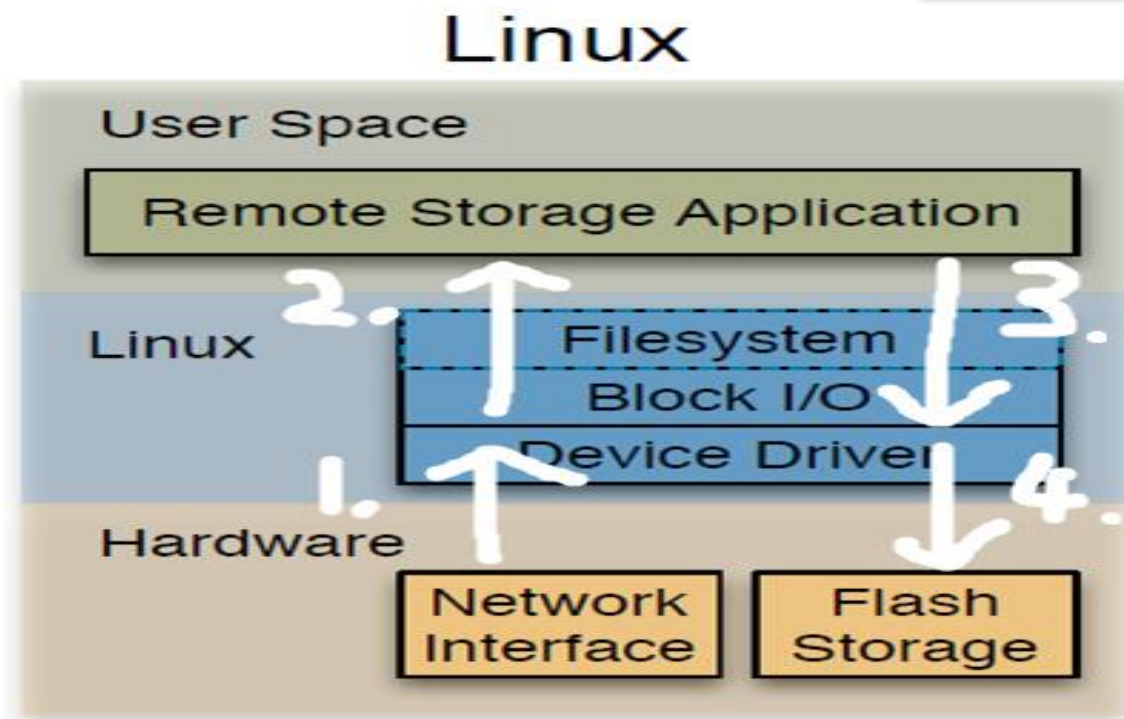
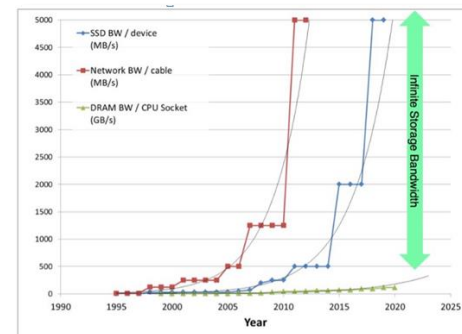
- Same data as last slide, but for the Log-impaired
- Storage Bandwidth is not literally infinite
- But the **ratio** of Network and Storage to CPU throughput is widening very quickly





异构硬件数据处理背景

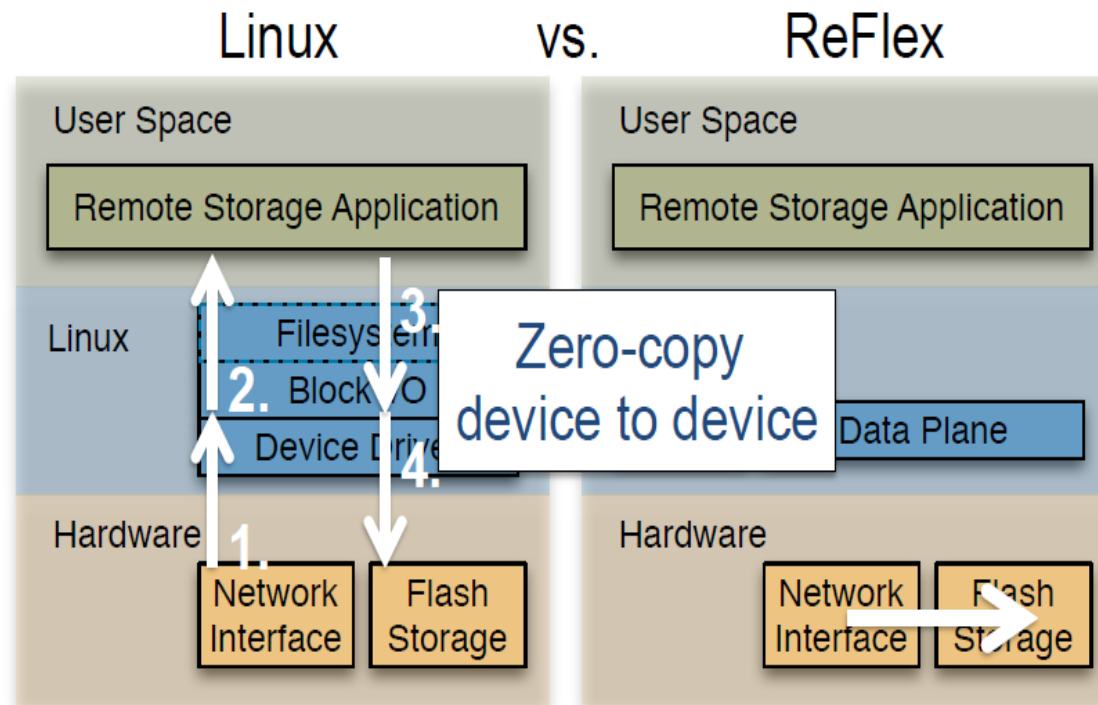
- 传统：中断式数据处理模式
 - 数据会经过多次拷贝





异构硬件数据处理背景

- 传统：中断式数据处理模式
 - 目标：数据零拷贝

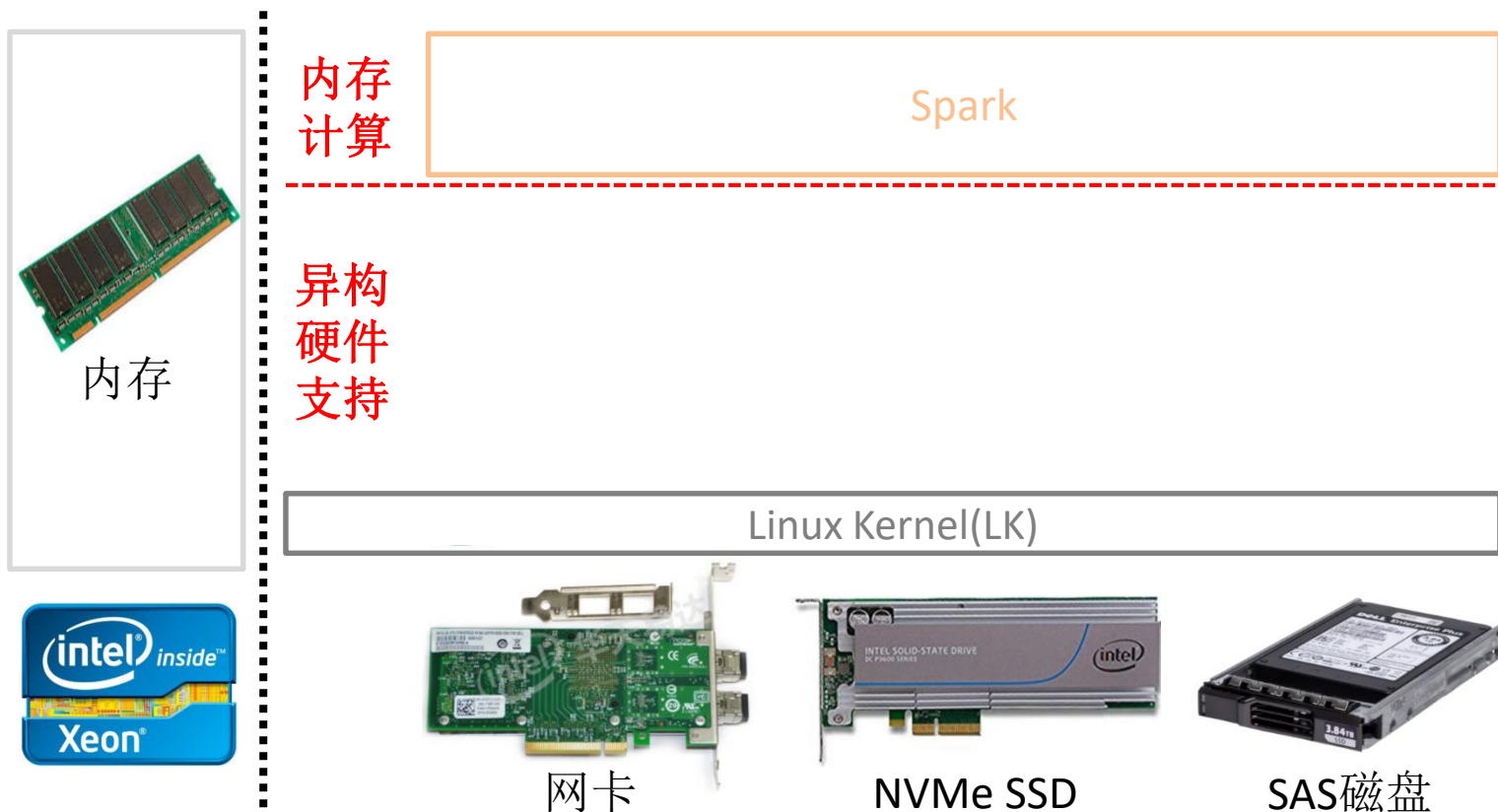




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

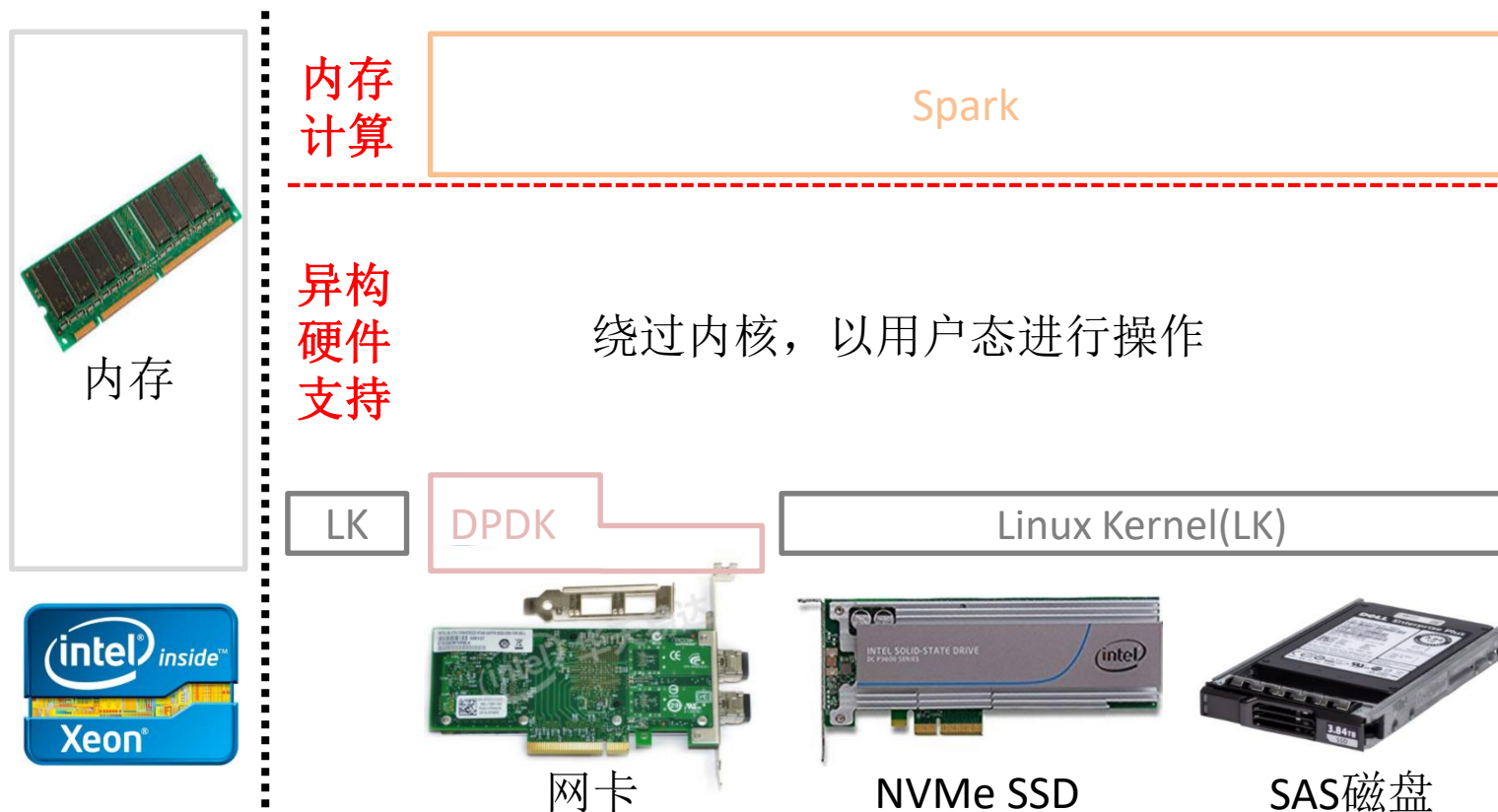




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

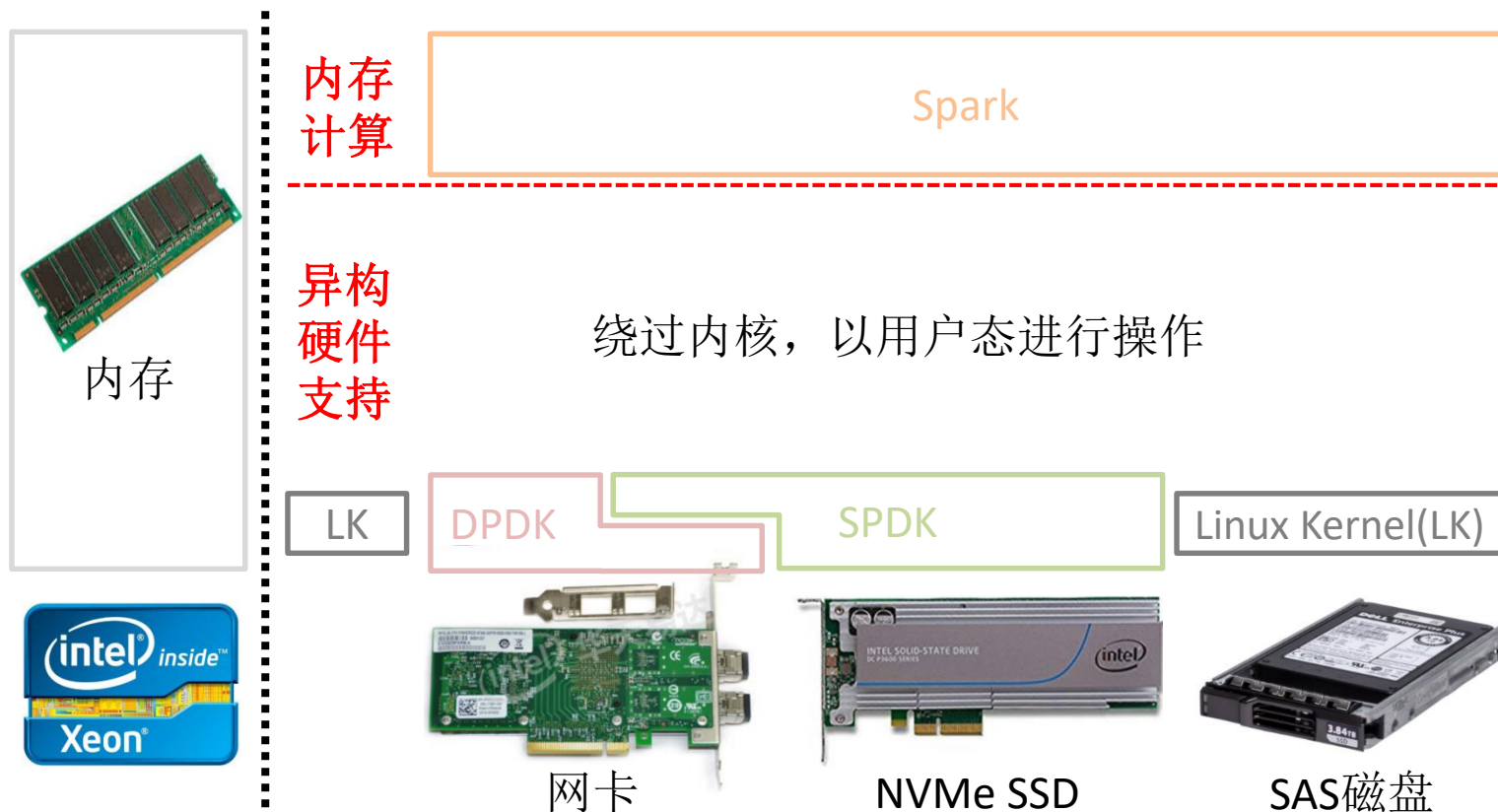




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

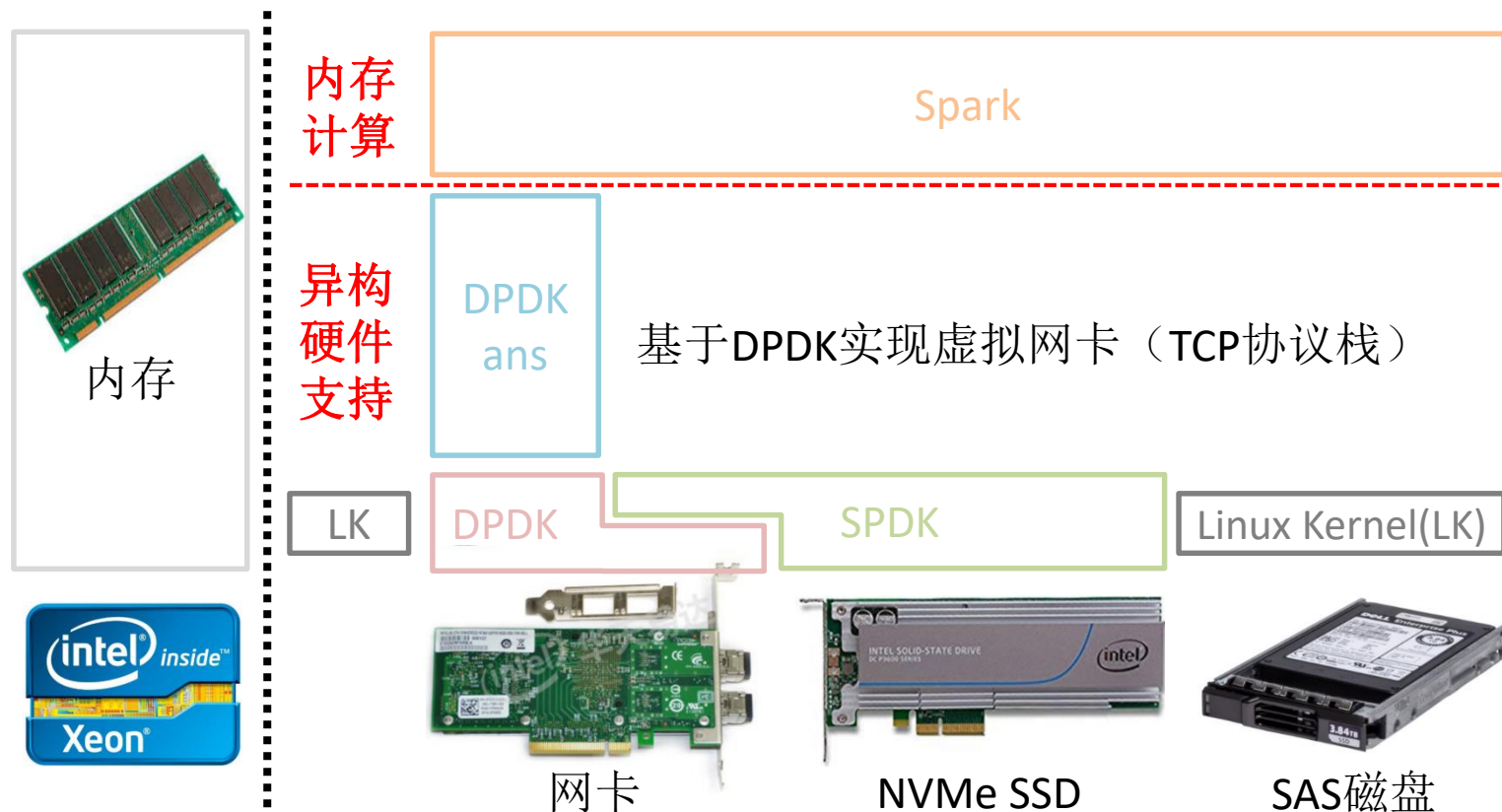




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

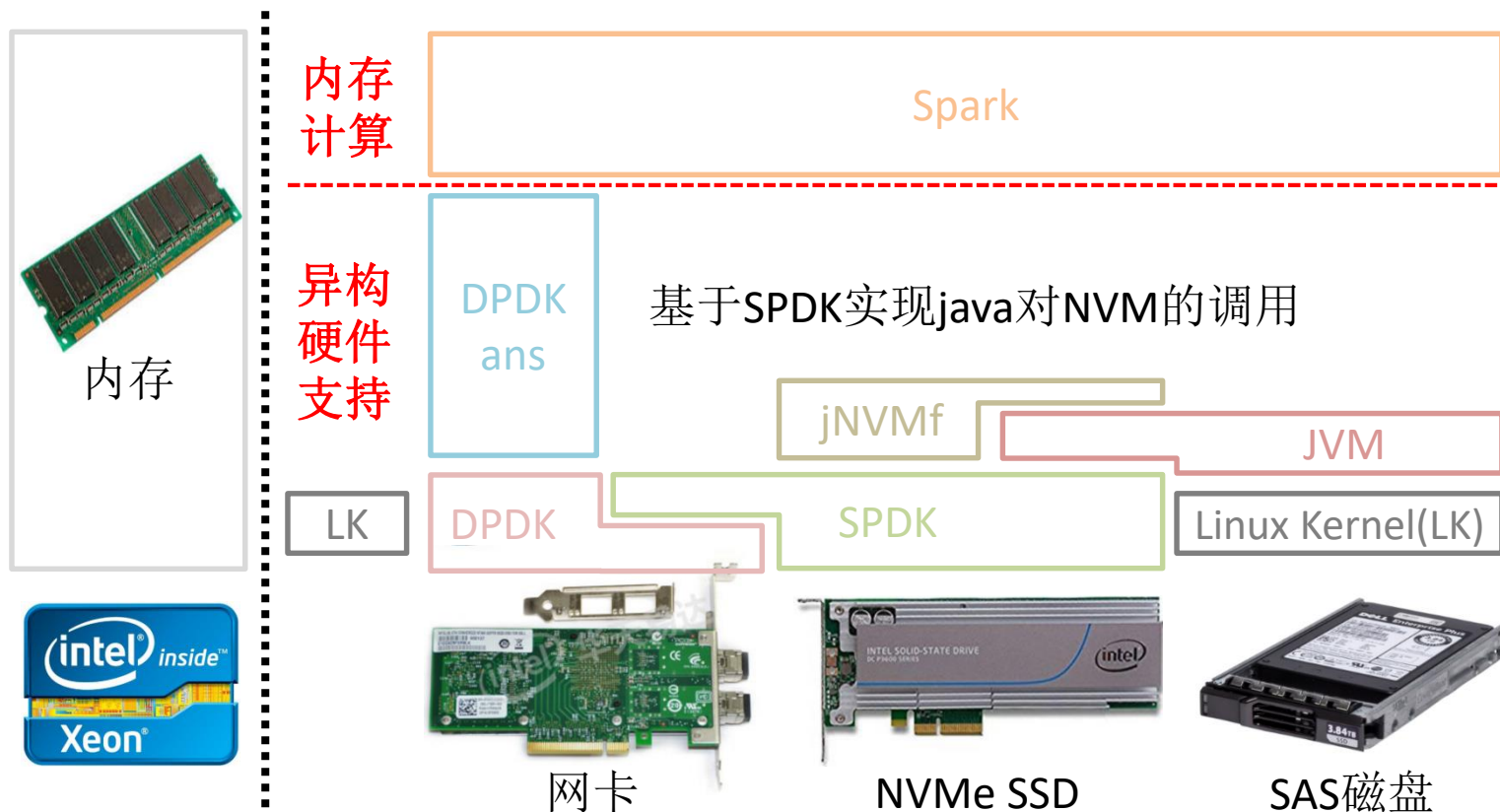




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

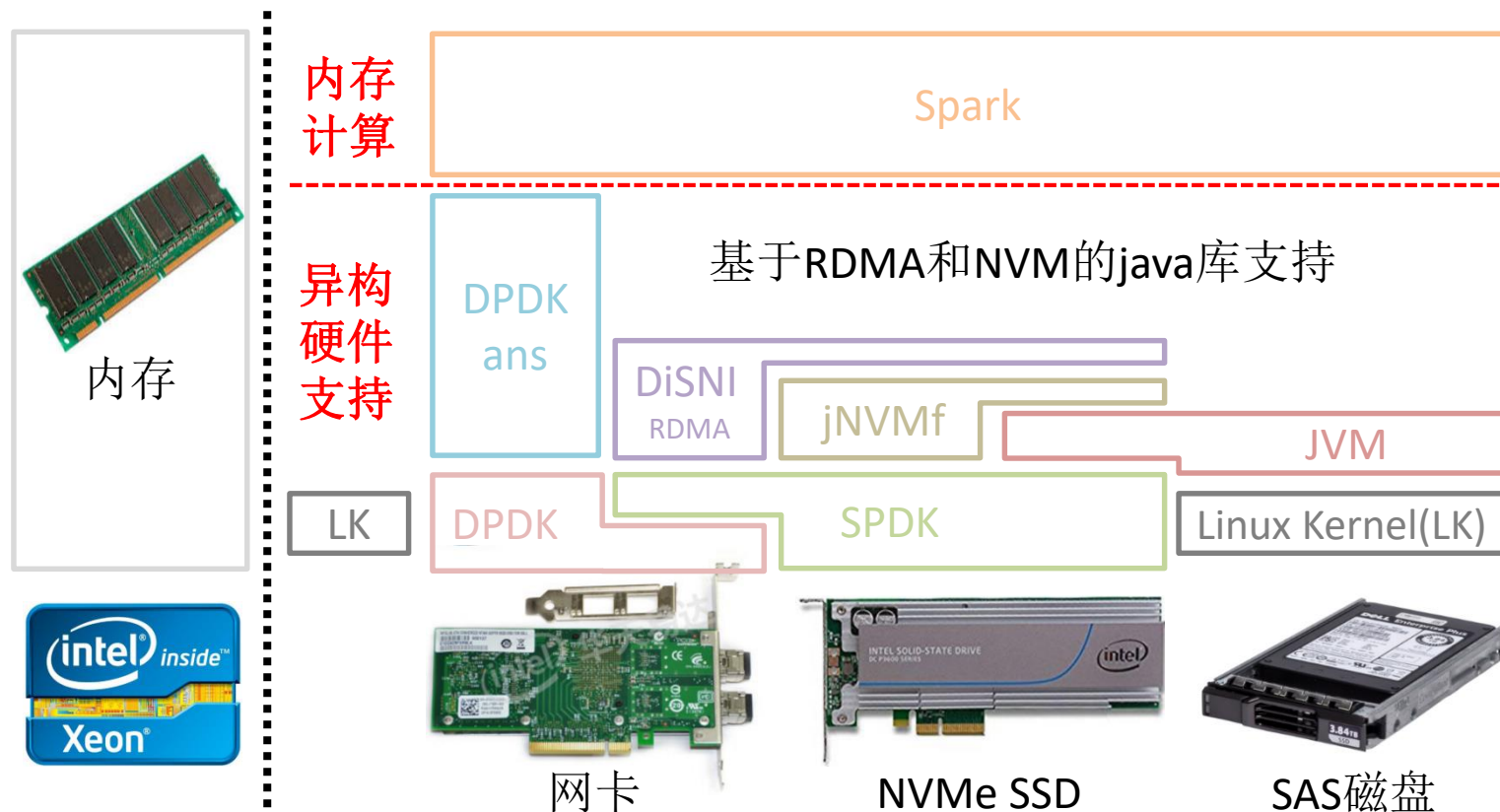




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

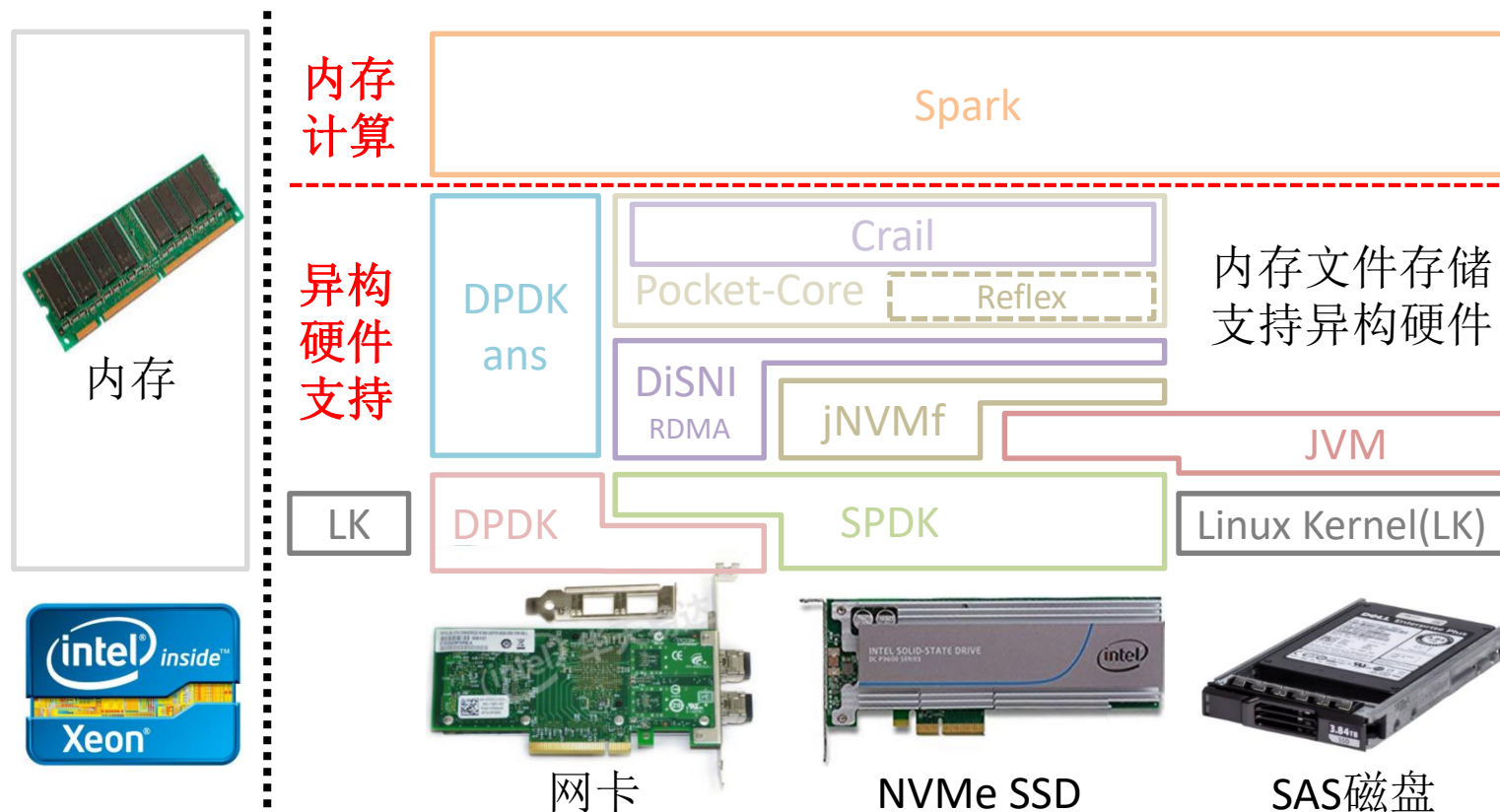




目标：零拷贝数据处理

异构硬件

- 基于异构硬件的数据处理

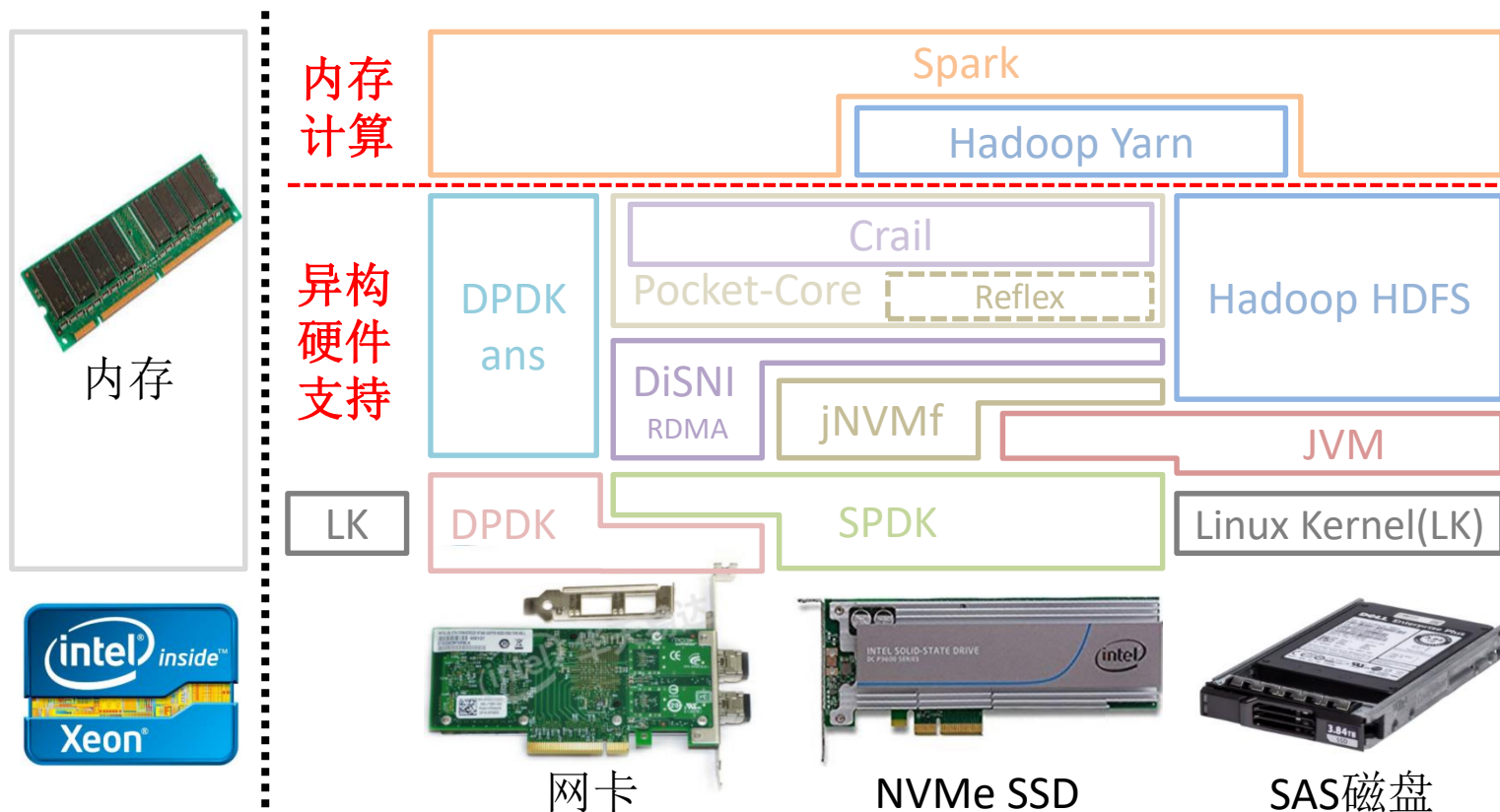




目标：零拷贝数据处理

异构硬件

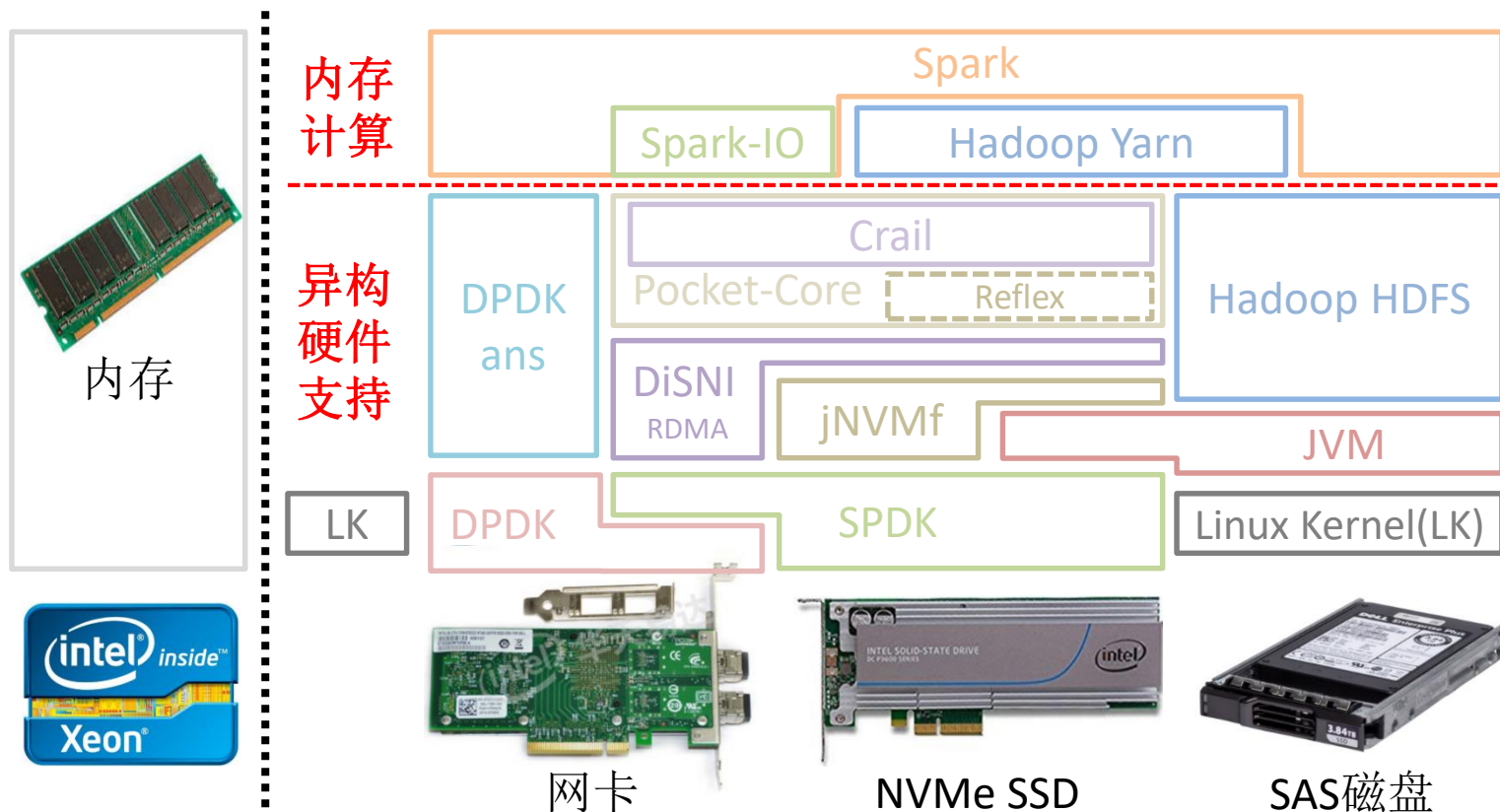
- 基于异构硬件的数据处理





零拷贝数据处理架构 异构硬件

- 基于异构硬件的数据处理





谢谢大家！

金熠波

2019年5月15日

逸夫楼C-115