

百度统一分布式计算框架Bigflow

百度基础架构部 张云聪

Bai du INF

Agenda

- 背景

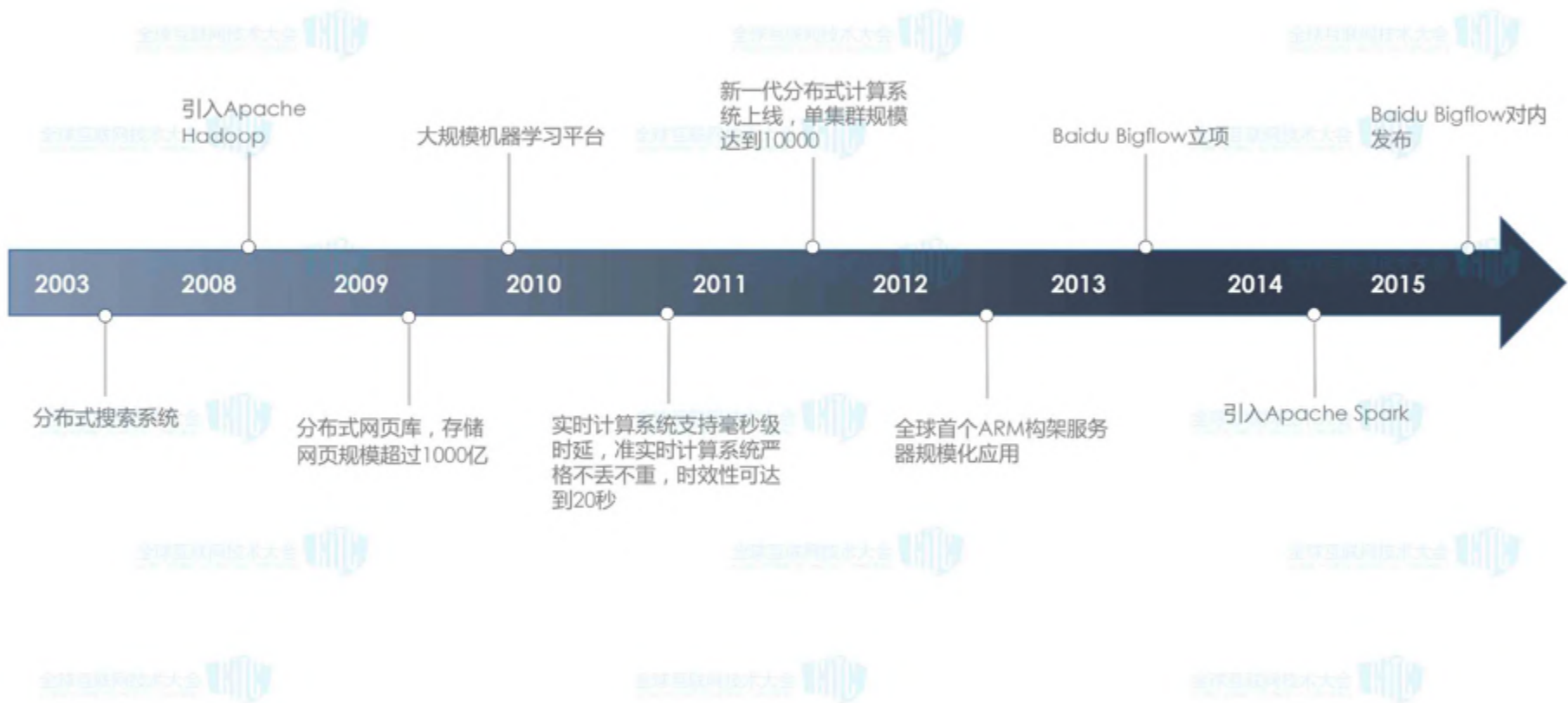
- 基本抽象/接口

- 统一离线/实时计算

- 优化

- 应用实践

背景



Batch

MapReduce
Abaci

Iterative

Spark

RealTime

DStream

ML

ELF

多个计算平台

开发慢 | 使用繁 | 维护难

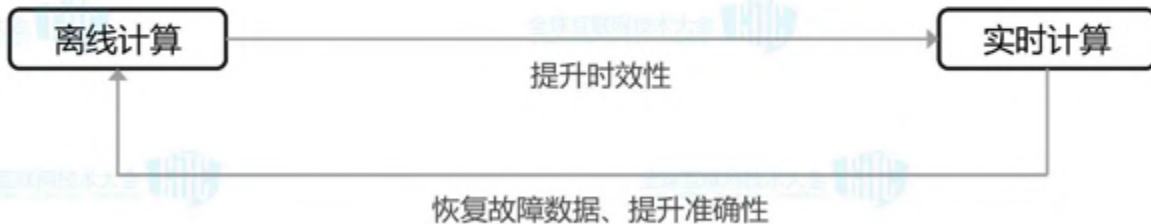
Mini-Batch

TaskManager
SparkStreaming

- 一套逻辑，不断重写：



- 一套逻辑，同时维护：



- Baidu MapReduce进化：

- 用户无感知：

流式Shuffle服务

Mapper/Reducer进程复用

Native Streaming

- 用户感知，推广难度大：

DAG拓扑支持

多输入多输出与Broadcast支持

并发动态调节

易用

高效

统一多引擎

简化模型

用户代码量短

方便本地调试

自动完成优化

自动调参

逻辑计划优化

运行时优化

高层次抽象

多引擎同一套接口

...

基本需求

细化

分布式统一 表示层

有了合理的表示层设计，

才能给前述所有需求提供大的发挥空间。

定位

指导思想

指导思想：使分布式程序尽可能像单机程序

MR => FlumeJava/Crunch => Spark/Flink/Beam



越来越函数式

函数式的分布式计算接口是社区趋势

指导思想 – 现有抽象的缺陷

指导思想：使分布式程序尽可能像单机程序。

细则一：不应暴露分布式的实现细节给用户。

细则二：所有的抽象应尽可能与单机语言中找到原型。

细则三：不应让用户去根据底层引擎特性来优化作业。

指导思想 – 现有抽象的缺陷

指导思想：使分布式程序尽可能像单机程序。

An RDD is a read-only, partitioned collection of records.

——Spark RDD论文

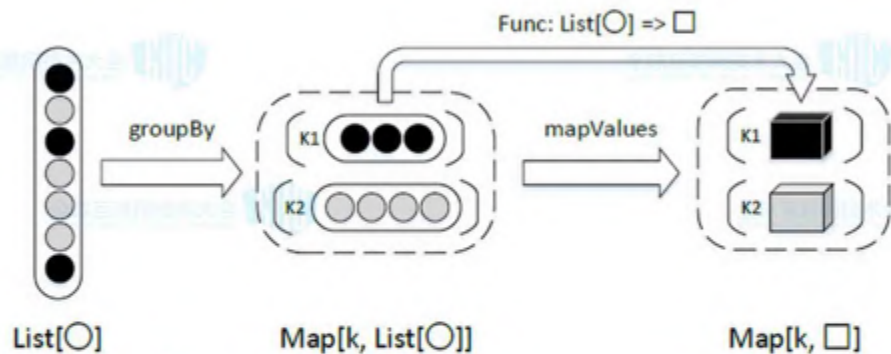


单机用户

什么是partition？

我为什么要关注数据是怎么分片的？

指导思想：使分布式程序尽可能像单机程序



单机：groupBy + reduce

现有分布式系统：reduceByKey

为什么我不能先分组，
然后再在每个组上聚合？



单机用户

指导思想：使分布式程序尽可能像单机程序

RDD中存储的元素为什么不能是一个RDD?



单机用户

```
hi, 咨询一下, 我现在想要解决的问题是在map里边做一个list转换为rdd, 这时候要用到JavaSparkContext sc.parallelize(list)的这种操作, 我把sc作为一个变量传到类里边,  
JavaPairRDD<String, JavaRDD<LabeledPoint>> groupRdd = groupPointRdd.mapValues(new Function<Iterable<LabeledPo  
{  
    private static final long serialVersionUID = 1L;  
  
    @Override  
    public JavaRDD<LabeledPoint> call(Iterable<LabeledPoint> v1) throws Exception {  
        // TODO Auto-generated method stub  
        List<LabeledPoint> list = new ArrayList<LabeledPoint>();  
        for (LabeledPoint v : v1) {  
            list.add(v);  
        }  
        JavaRDD<LabeledPoint> labeledPointRdd = sc.parallelize(list);  
        return labeledPointRdd;  
    }  
});
```

为transient, 但我标的话这时候就会报sc空指针, 各位老师有啥这方面的解决方案吗? @张三聊

不支持嵌套, 每个组的数据就不能被表达为RDD。
则现有分布式算法无法复用于每个分组。

指导思想：使分布式程序尽可能像单机程序

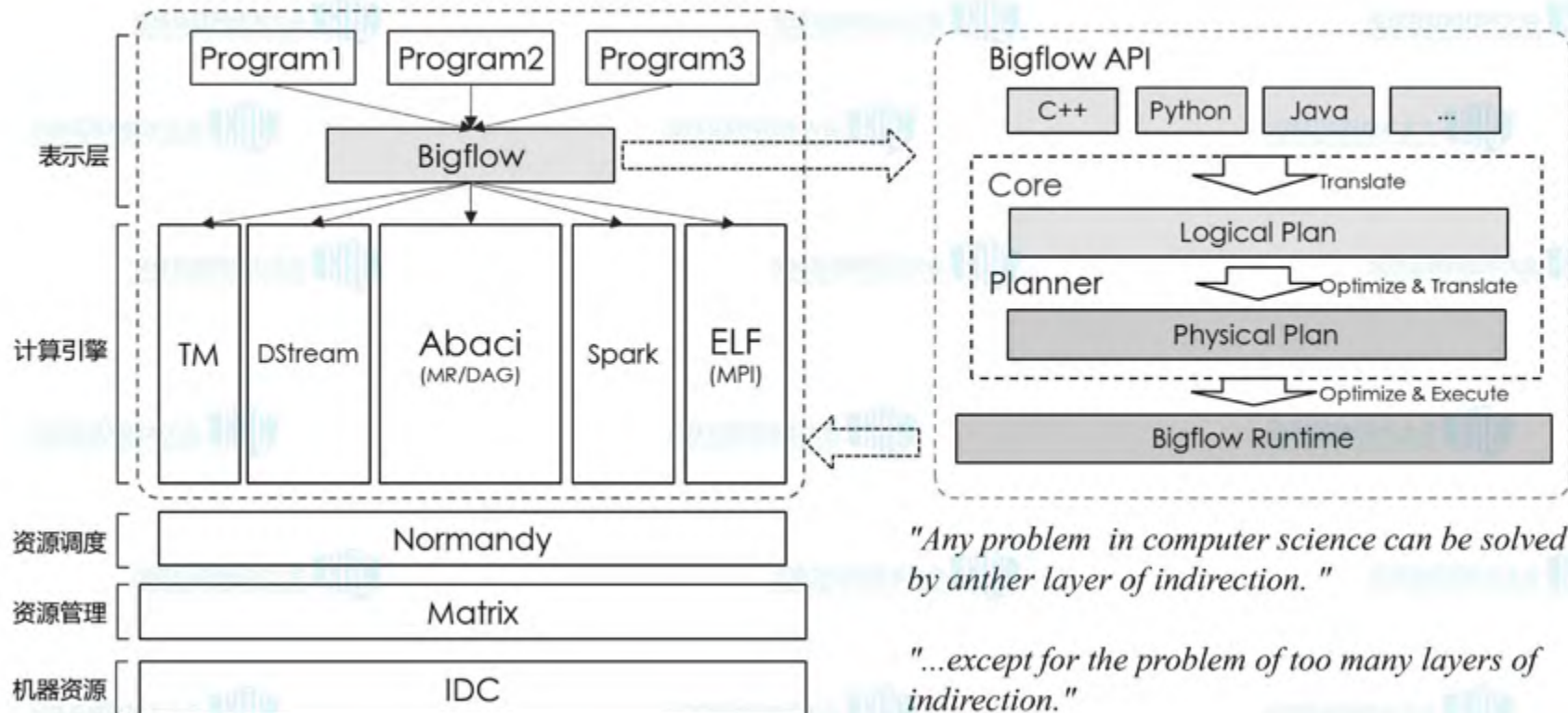
Bigflow

引入单机常见的
嵌套数据集
概念

删除单机没有的
Partition等
概念

去掉冗余的
reduceByKey等
接口

合理的接口设计
带来更大的优化
空间



Agenda

- 背景
- 基本抽象/接口
- 统一离线/实时计算
- 优化
- 应用实践



- Pipeline: 入口 / 分布式作业的抽象
 - MRPipeline/SparkPipeline...
- Dataset: 分布式数据集抽象
 - 有序/无序 by **Sort**
 - 单元素/任意数量元素 by **Aggregations**
 - 有穷/无穷 from **Source/Window**
 - 有/无Schema
 - 扁平/嵌套 by **Groupings**
- Transformations: 变换
 - 全部惰式计算
 - **SideInput**

```
p = base.Pipeline.create(ENGINE)
```

```
words = p.read(input.TextFile('/a/b/c')) \  
        .flat_map(lambda line: line.split('\t'))
```

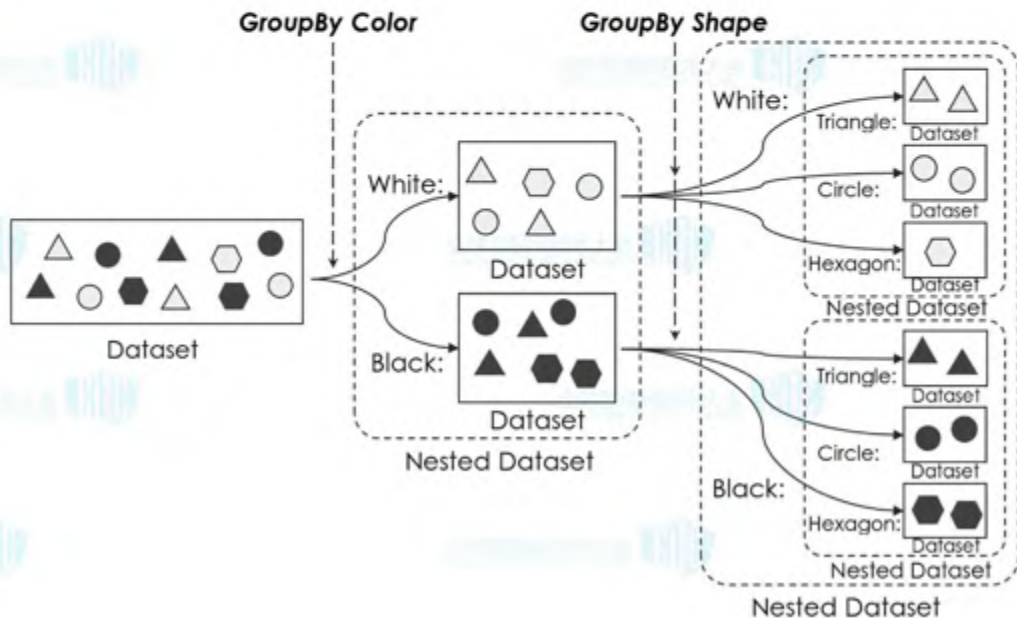
```
per_word_cnt = words \  
               .group_by(lambda word: word) \  
               .apply_values(count) \  
               .flatten()
```

```
p.write(word_cnt, output.TextFile('/a/b/e'))  
p.run()
```


- DataSet嵌套分组

- `DataSet<Color, DataSet<Element>> grouped = dataset.apply(GroupBy, getColor)`

- `DataSet<Color, DataSet<Shape, DataSet<Element>>> nested = grouped.apply_values(GroupBy, getShape)`



```
p = base.Pipeline.create(ENGINE)

visitors = p.read(input.TextFile('/a/b/c')) \
    .flat_map(lambda line: line.split('\t'))

def cnt_distinct(visitors):
    """ @param visitors: DataSet<String> """
    return visitors \
        .distinct() \
        .count()

num = visitors.apply(cnt_distinct)

p.write(num, output.TextFile('/a/b/d'))
p.run()
```

```
total_visitors = ...

def uv_counting(total_visitors):
    """ @param page_visitors: DataSet<Pair<String, String>> """
    return total_visitors \
        .group_by(lambda (visitors, page): page) \
        .apply_values(cnt_distinct) \
        .flatten()

pv_result = total_visitors.apply(uv_counting)

p.write(num, output.TextFile('/a/b/d'))
p.run()
```

Transformation	UDF	Input => output
map(p, f)	$f: (V) \Rightarrow T$	$[V] \Rightarrow [T]$
filter(p, f)	$f: (V) \Rightarrow \text{bool}$	$[V] \Rightarrow [V]$
flat_map(p, f)	$f: (V) \Rightarrow \text{Iterable}[T]$	$[V] \Rightarrow [T] / V \Rightarrow [T]$
group_by_key(p)	-	$[(K, V)] \Rightarrow \{K: [V]\}$
group_by(p, f)	$f: (V) \Rightarrow K$	$[V] \Rightarrow \{K: [V]\}$
flatten(p)	-	$\{K: [V]\} \Rightarrow [(K, V)]$
reduce(p, f)	$f: (V, V) \Rightarrow V$	$[V] \Rightarrow V$
aggregate(p, v, f1, f2)	$v: T \text{ or } () \Rightarrow T$ $f1: (T, V) \Rightarrow T$ $f2: (T, T) \Rightarrow T$	$[V] \Rightarrow T$
join(p1, p2)	-	$[(K, V1)], [(K, V2)] \Rightarrow [(K, (V1, V2))]$
cogroup(p1, p2)	-	$[(K, V1)], [(K, V2)] \Rightarrow \{K: ([V1], [V2])\}$
cartesian(p1, p2)	-	$[V1], [V2] \Rightarrow [(V1, V2)]$
union(p1, p2)	-	$[V], [V] \Rightarrow [V]$
sum(p)	-	$[V] \Rightarrow V$
count(p)	-	$[V] \Rightarrow \text{Long}$
take(p, num)	-	$[V] \Rightarrow [V]$
sort(p, f)	$f: f(V) \Rightarrow K$	$[V] \Rightarrow [V]$

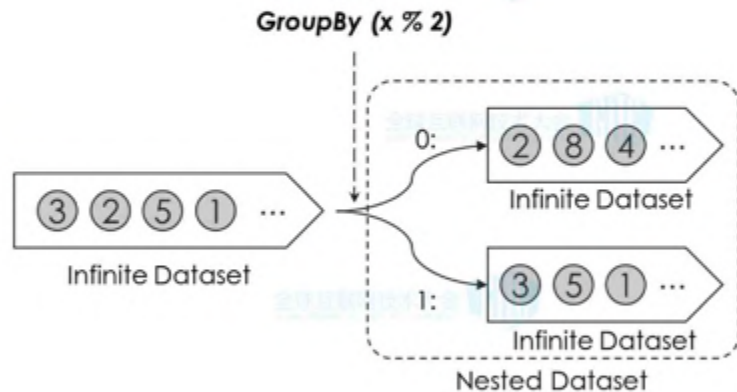
Agenda

- 背景
- 基本抽象/接口
- 统一离线/实时计算
- 优化
- 应用实践



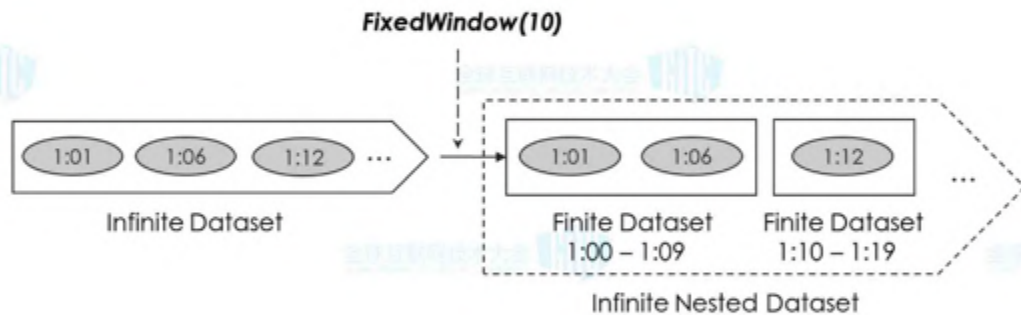
抽象—Nestable Dataset

- 无穷DataSet分组



➤ `DataSet<Odeivity, DataSet<Int>> grouped = dataset.apply(GroupBy, modByTwo)`

- 无穷DataSet划分窗口



➤ `DataSet<Window, DataSet<Log>> windowed = dataset.apply(WindowInto, FIXED_WINDOW, getTime)`

```
total_visitors = ...
```

```
def uv_counting(total_visitors):
```

```
    return total_visitors \
```

```
        .group_by(lambda (visitors, page): page) \
```

```
        .apply_values(cnt_distinct)
```

```
        .flatten()
```

```
pv_result = total_visitors.apply(uv_counting)
```

```
p.write(num, output.TextFile('/a/b/d'))
```

```
p.run()
```

```
def uv_counting(total_visitors):
```

```
    return total_visitors \
```

```
        .group_by(lambda (visitors, page): page) \
```

```
        .apply_values(cnt_distinct)
```

```
        .flatten()
```

```
raw_words = p.read(input.Kafka(xxx)) \
```

```
    .flat_map(lambda line: line.split())
```

```
num = words \
```

```
    .window_into(time_extractor, FIXED_WINDOW) \
```

```
    .apply_values(uv_counting) \
```

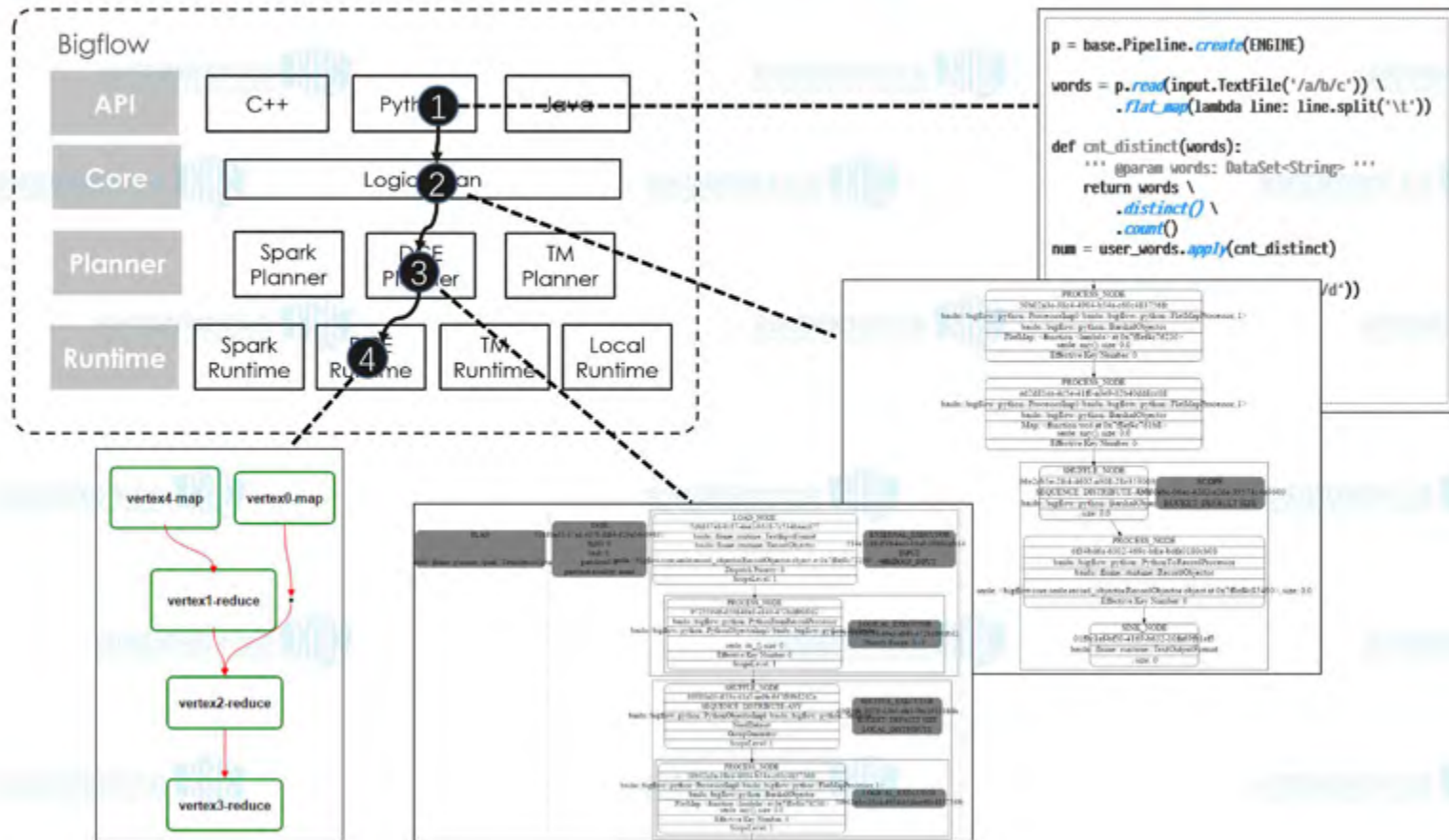
```
    .flatten()
```

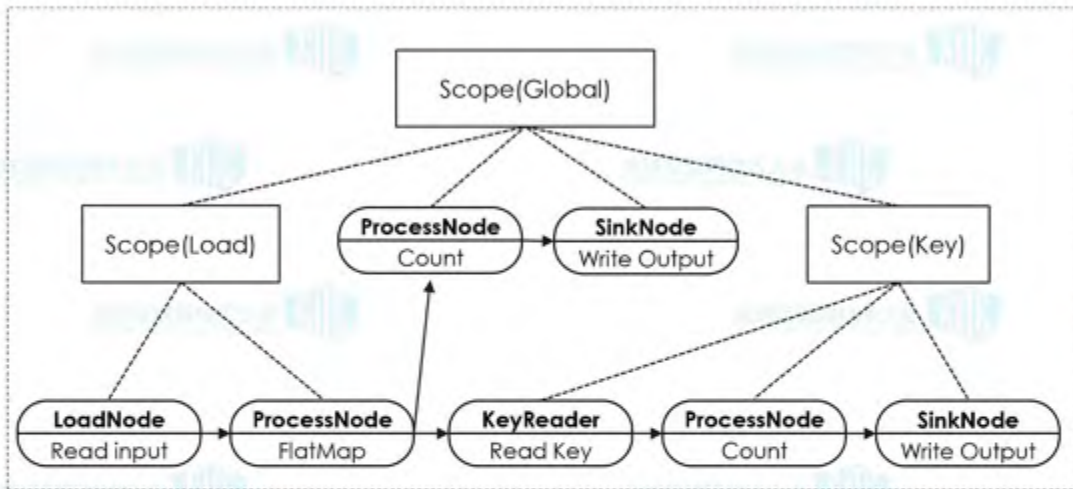
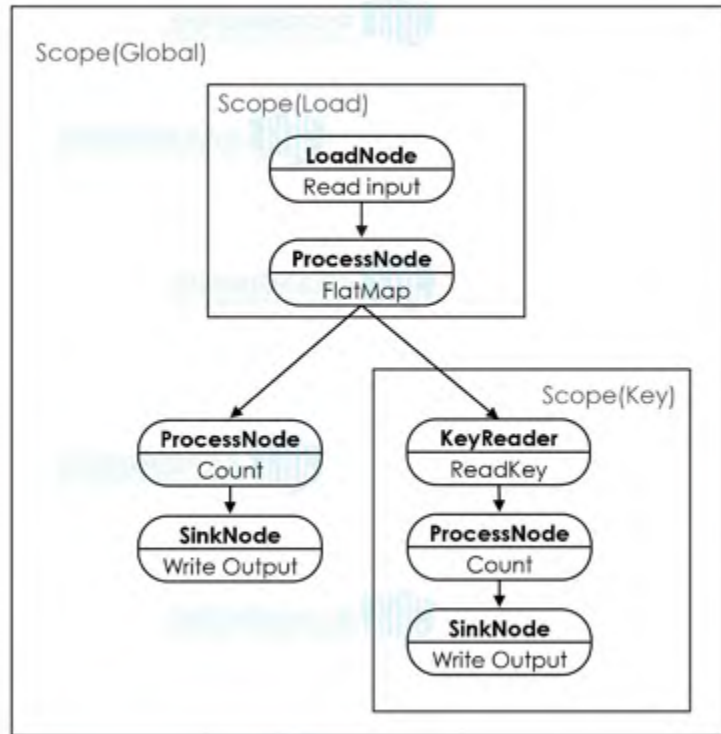
```
p.write(num, ...)
```

```
p.run()
```

Agenda

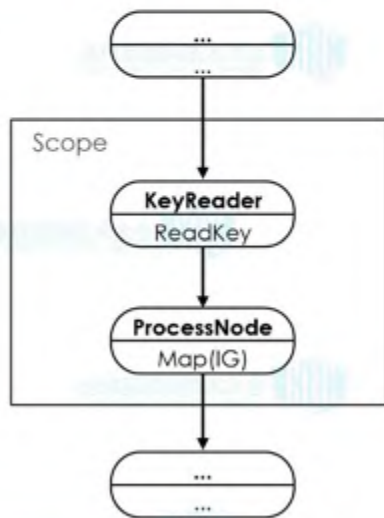
- 背景
- 基本抽象/接口
- 统一离线/实时计算
- 优化
- 应用实践



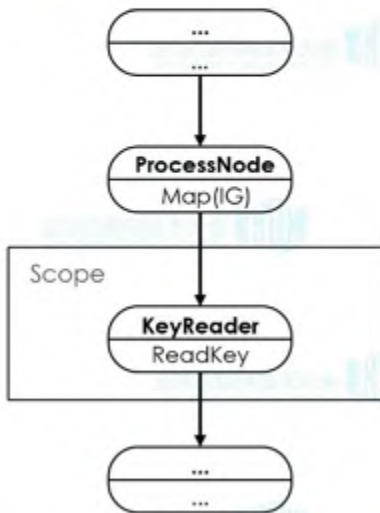


- Node: “计算” (DAG)
- Scope: “分布式” (Tree)

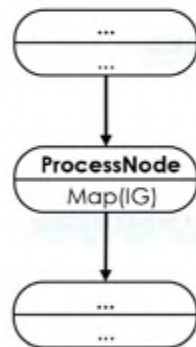
```
dataset \  
.group_by(_) \  
.apply_values(map, fn) \  
.flatten()
```



```
dataset \  
.map(fn) \  
.group_by(_) \  
.flatten()
```



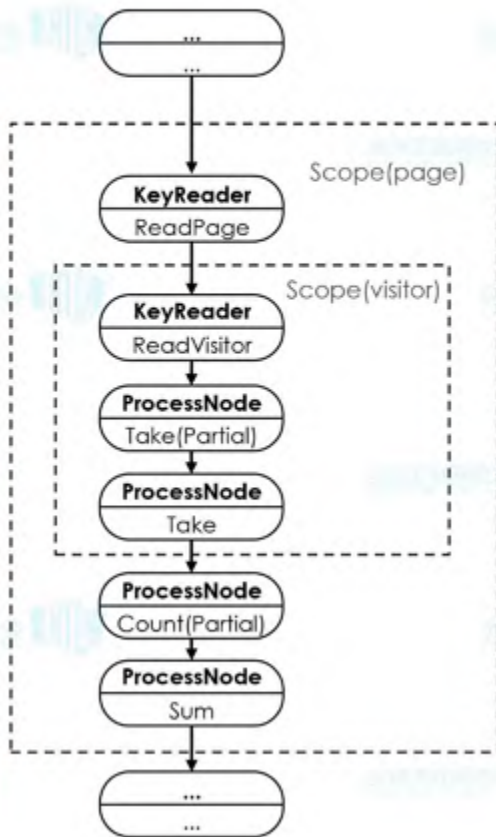
```
dataset \  
.map(fn)
```

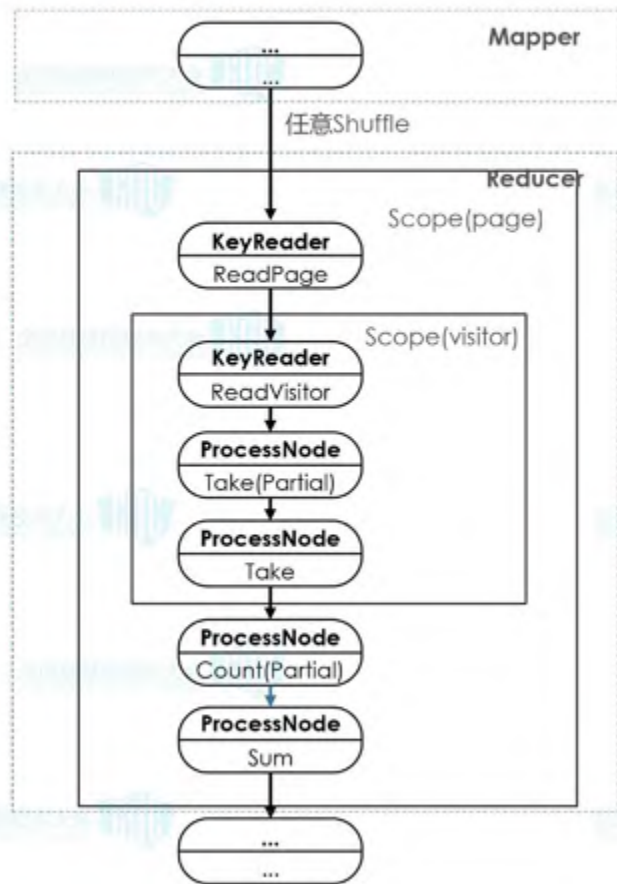


```
dataset \  
.group_by(lambda (visitors, page) : page) \  
.apply_values(cnt_distinct) \  
.flatten()
```

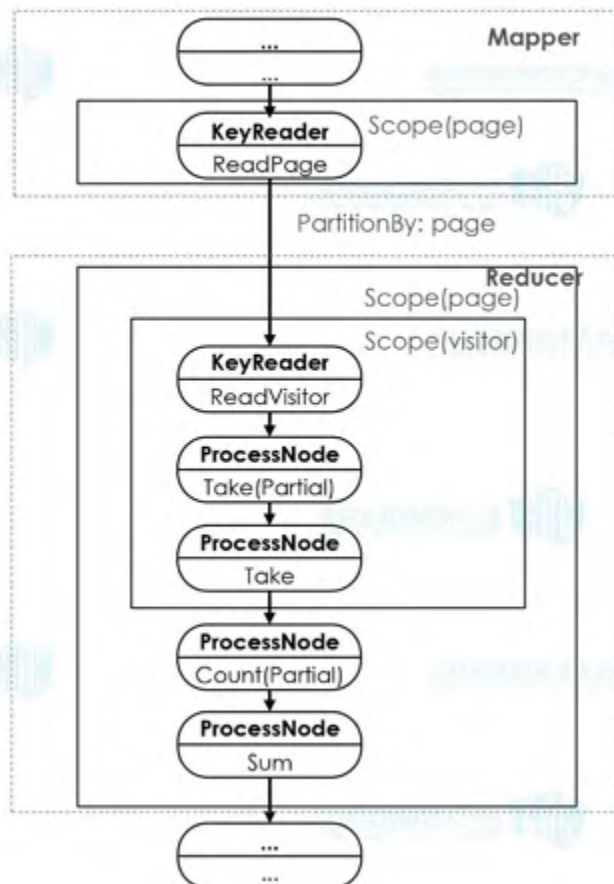


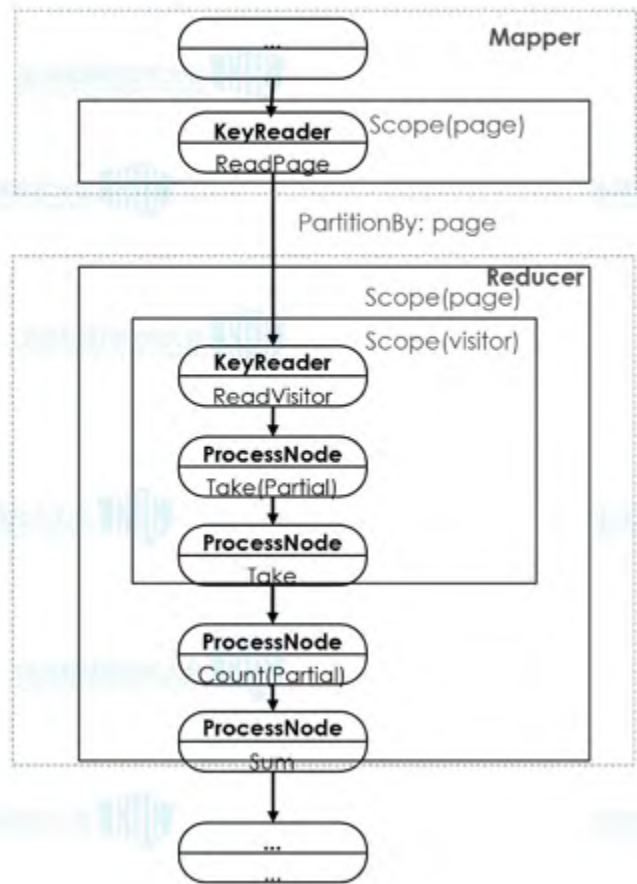
```
dataset \  
.group_by(lambda (visitors, page): page) \  
.apply_values(lambda visitors: visitors \  
.group_by(visitor) \  
.apply_values(take, 1) \  
.flatten() \  
.count()) \  
.flatten()
```



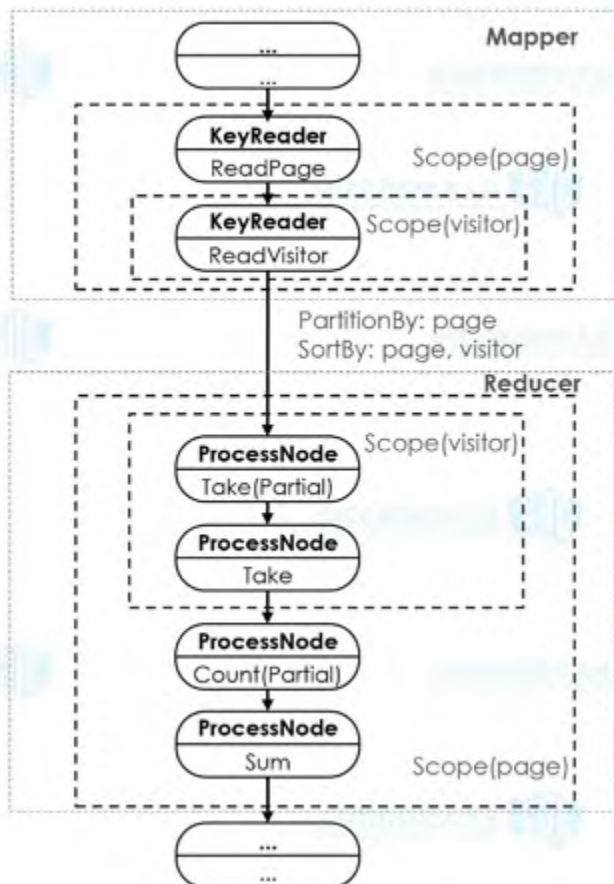


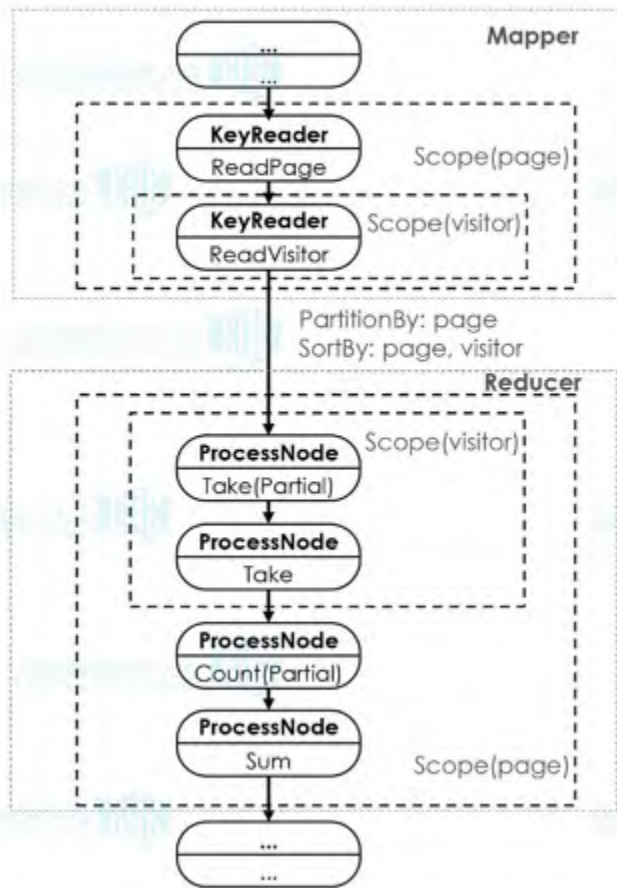
ShuffleScope
前移



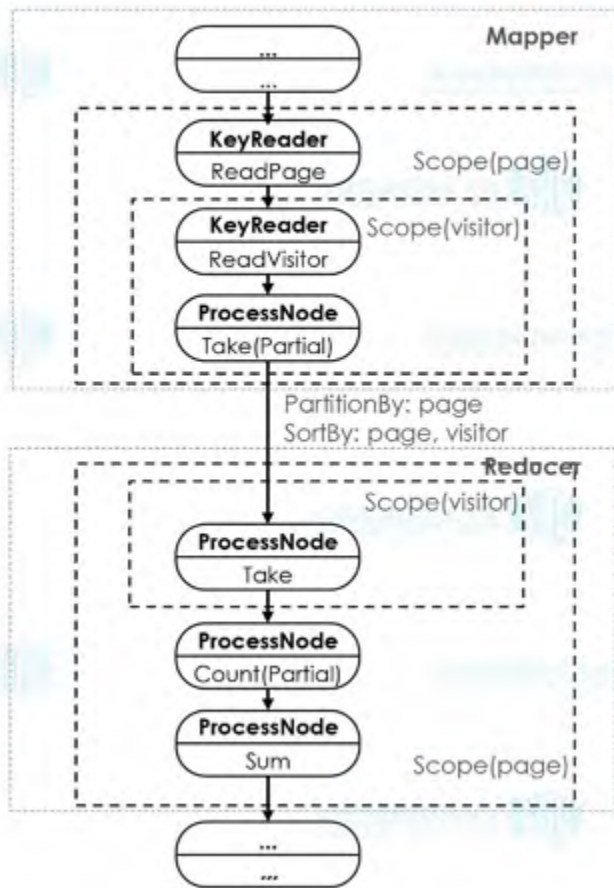


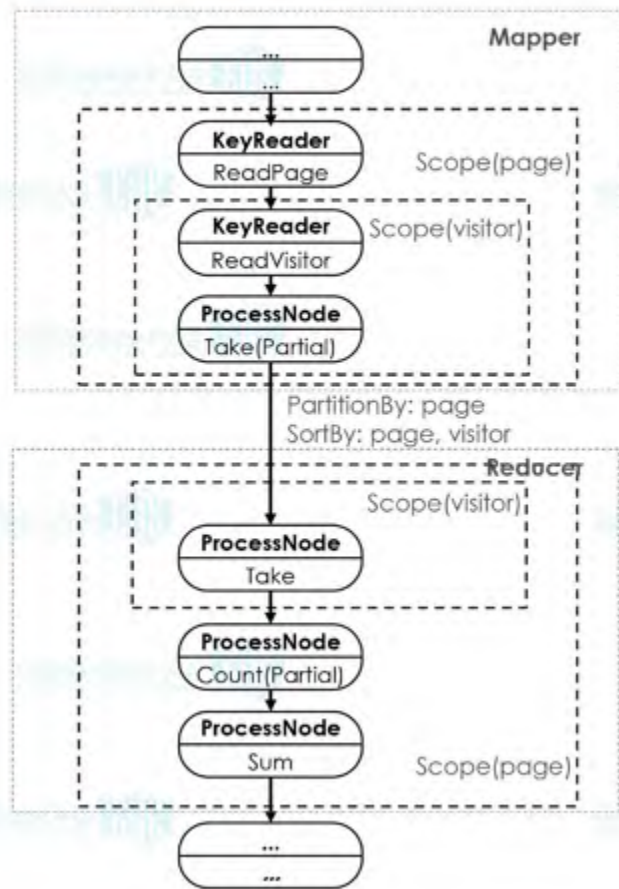
ShuffleScope
前移



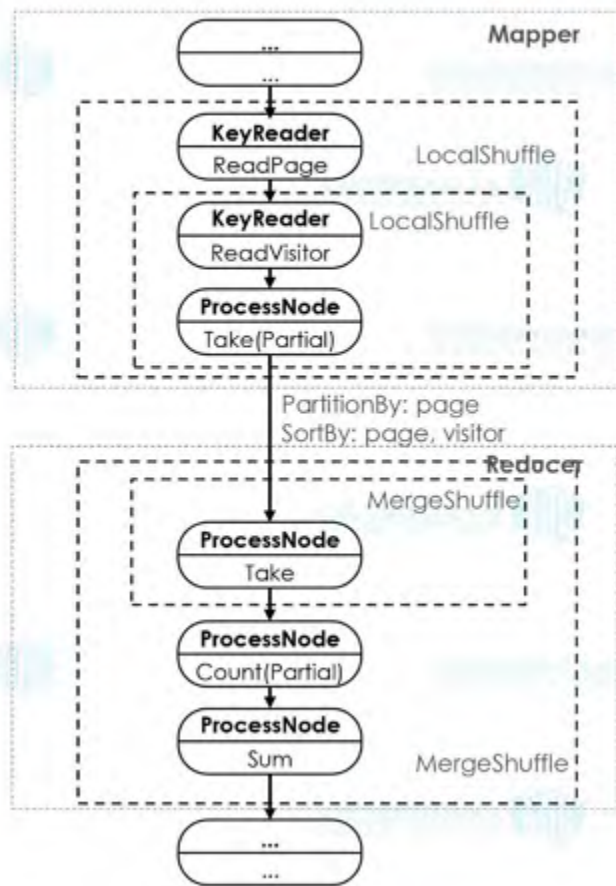


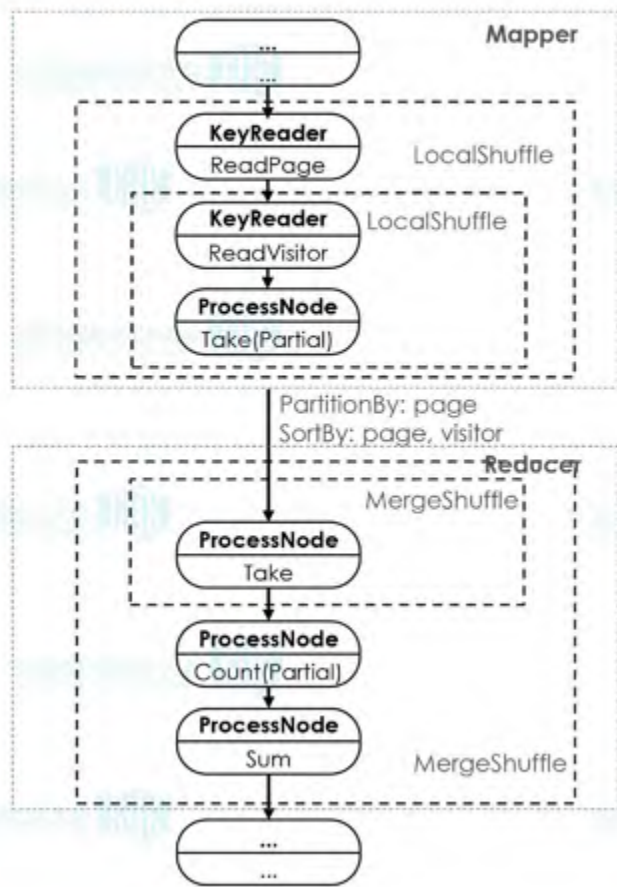
Partial算子前移



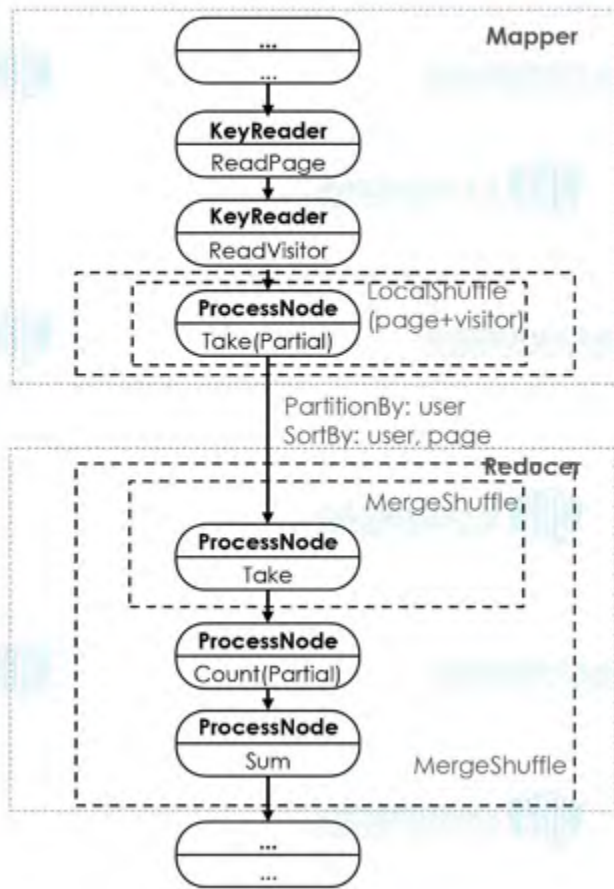


Shuffle类型
判定





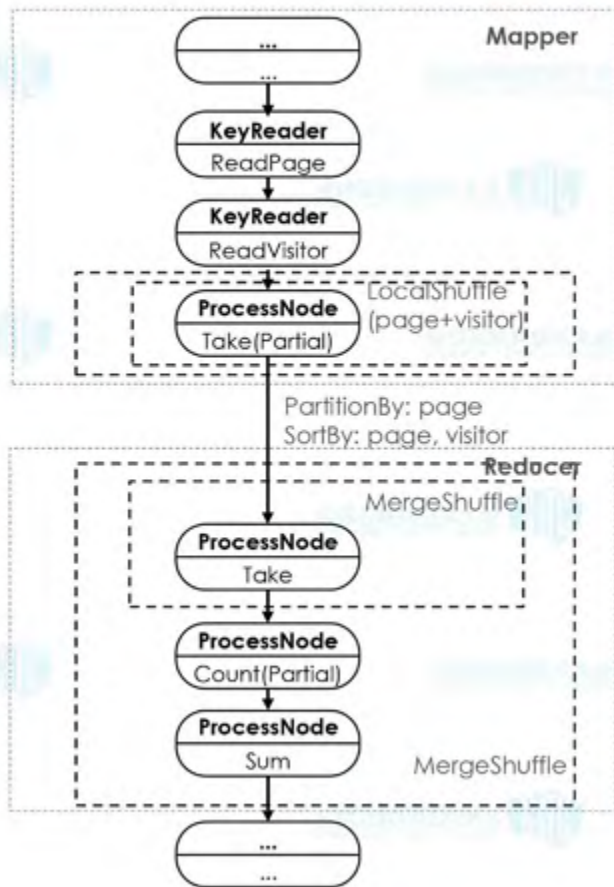
LocalShuffle
合并



```
dataset \
  .group_by(lambda (visitors, page) : page) \
  .apply_values(cnt_distinct) \
  .flatten()
```

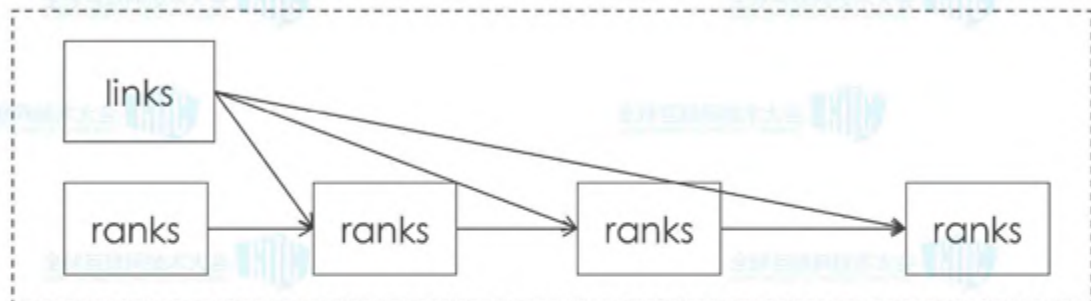


```
dataset \
  .group_by(lambda (visitors, page): page) \
  .apply_values(lambda visitors: visitors \
    .group_by(visitor) \
    .apply_values(take, 1) \
    .flatten() \
    .count()) \
  .flatten()
```

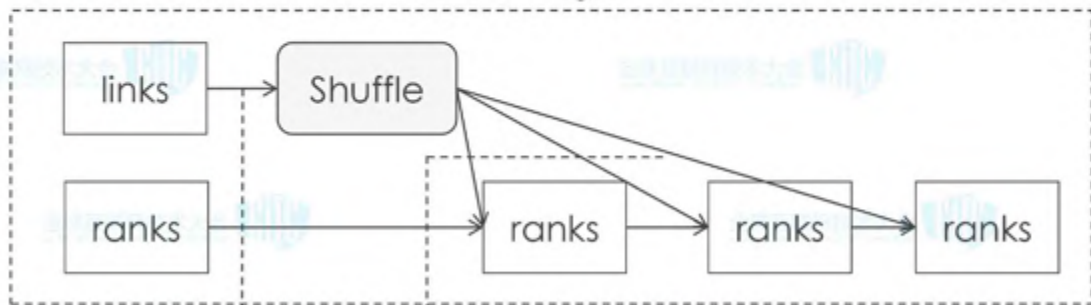


优化-示例(结合引擎)

- 示例: PageRank



同数据Shuffle合并



一次写

多次读

Agenda

- 背景
- 基本抽象/接口
- 统一离线/实时计算
- 优化
- 应用实践

应用实践-收益

日运行作业数3K+

日处理数据量3P+

在百度内部应用于上百产品线

数百活跃用户

Bigflow作业占比持续稳步提升



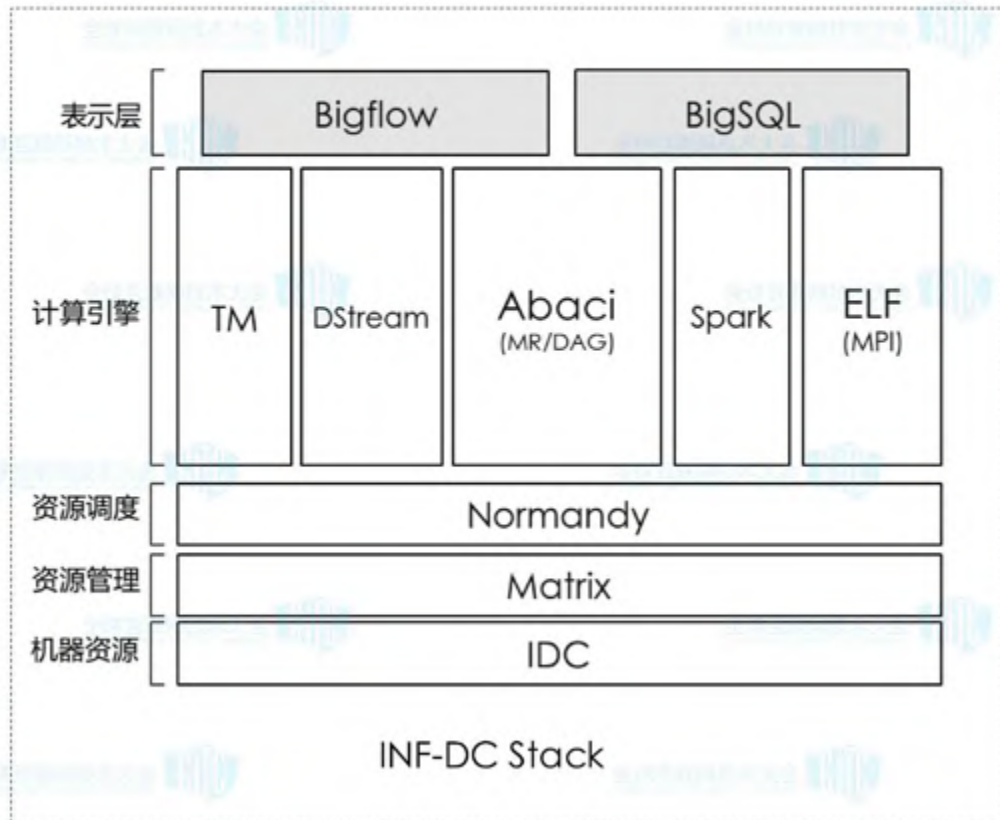
- 性能提升
 - C++ vs Hadoop Streaming
 - 算子前移: filter/select/初段aggregate
 - Python解释器复用, 相同字典只加载一次
- 多引擎支持
 - Hadoop MapReduce -> DAG MapReduce -> Spark
- 实际性能对比结果
 - 单个作业平均20%, DAG作业提升50%

BIGFLOW改写作业性能对比





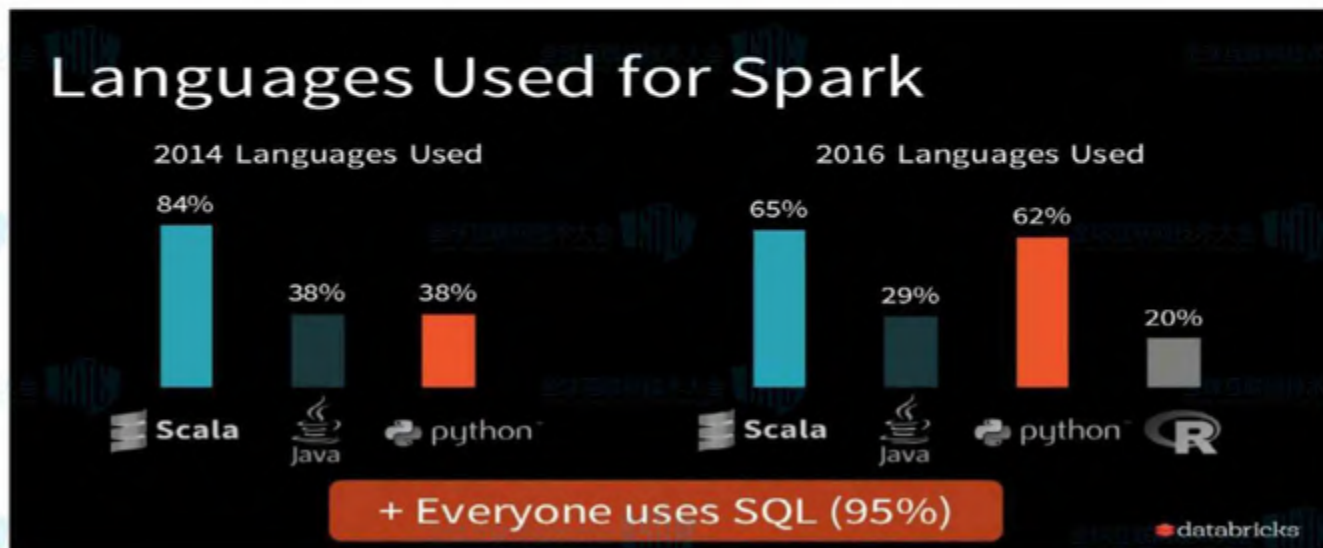
- 100+ 业务线, 300+ 用户
- 旧作业改写:
 - 30%代码量减少
 - 50%性能提升



目标：

2017 Q3 开源Bigflow on Spark

Python接口性能可大大高于PySpark



THANK YOU
Baidu Infrastructure