



# DISQL 2.0 : 百度海量数据分析语言

陈晓鸣  
资深工程师  
百度基础架构部  
[@陈晓鸣在百度](#)

[chenxiaoming@baidu.com](mailto:chenxiaoming@baidu.com)





发展历程

一个例子

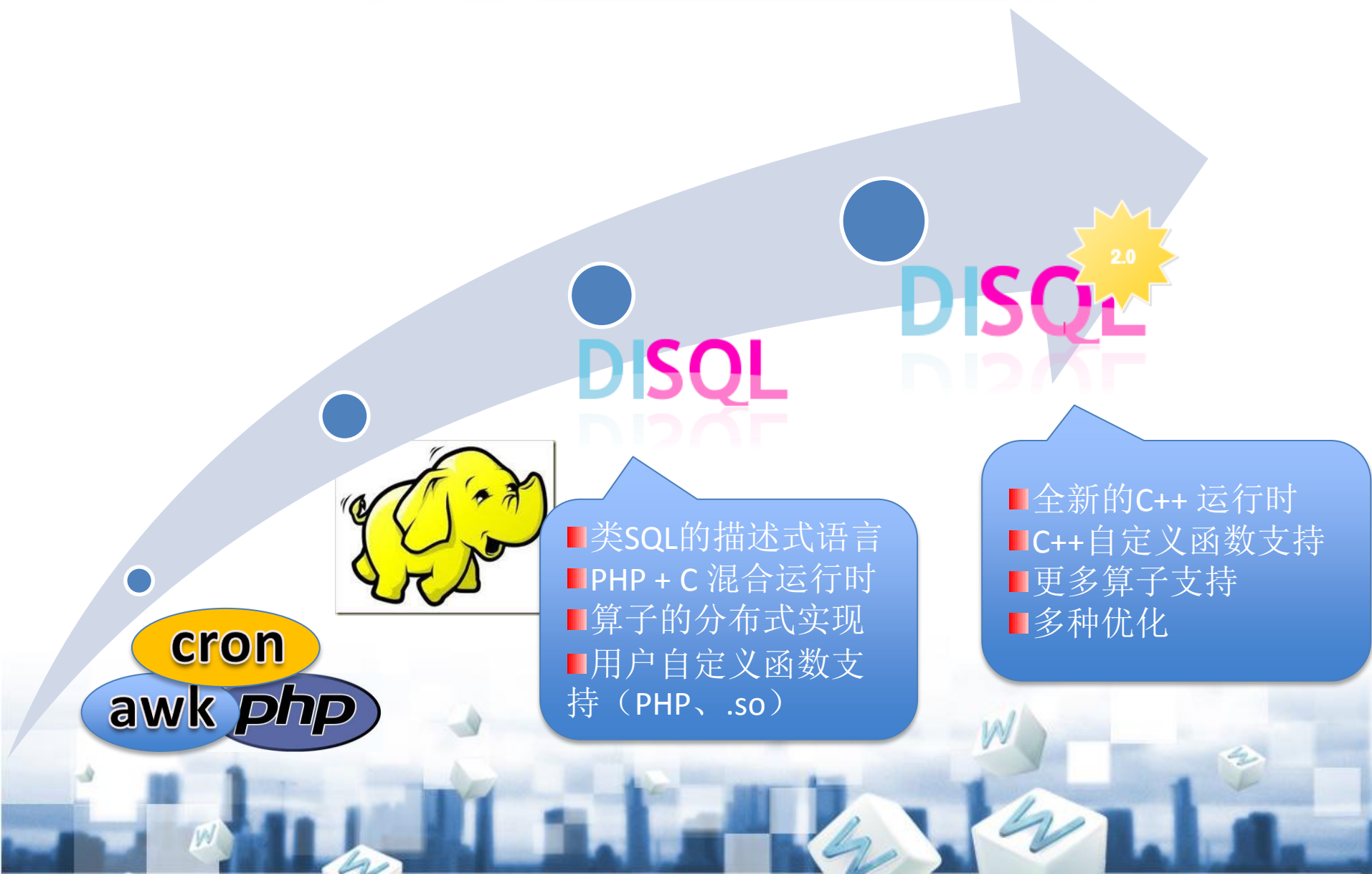
前端处理

中间语言翻译

运行时

总结与问答





发展历程



一个例子

前端处理

中间语言翻译

运行时

总结与问答

# 例：新闻站点访问量和广告量统计

## ■ 执行步骤

- 读取日志数据
- 选取出\_url、\_Res(广告数)两列
- 编写一个函数，从\_url中抽取出\_Site
- 用正则表达式过滤出新闻站点的数据
- 按站点分组，每组做两件事：
  - 计算访问量
  - 将广告数求和
- 输出数据，每行是一个JSON数据



## ■ 读取日志数据

■ 选取 \_Url、\_Res(广告数)两列

## ■ 编写一个函数，从 Url 中抽取出 Site

## ■ 用正则表达式过滤出新闻站点的数据

■ 按站点分组，每组做两件事：

## ■ 计算访问量

## ■ 将广告数求和

## 输出数据，每行是一个JSON数据

前端语  
言处理

中间语  
言翻译

运行时



发展历程

一个例子



前端处理

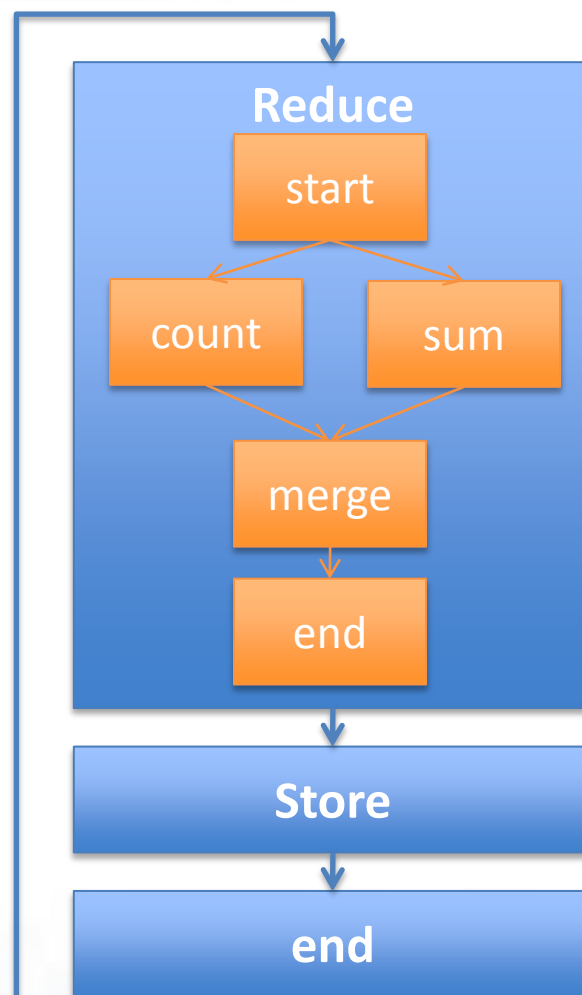
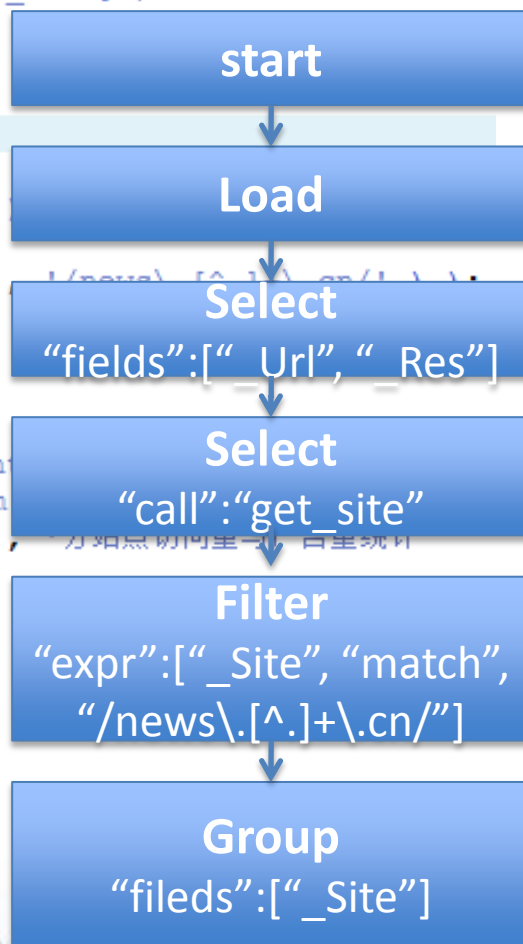
中间语言翻译

运行时

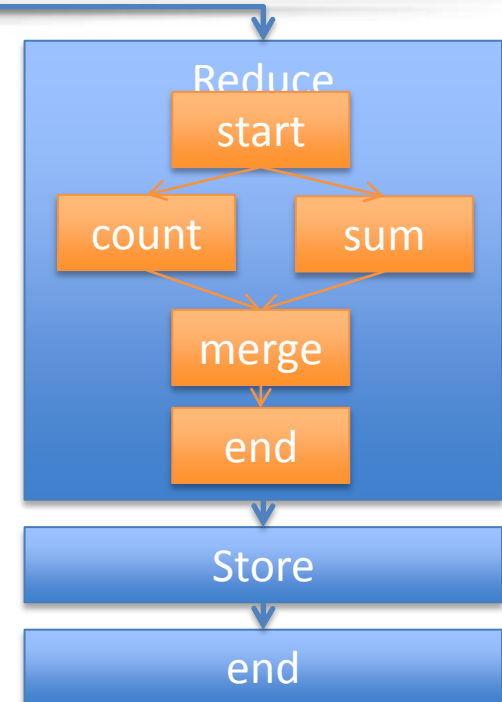
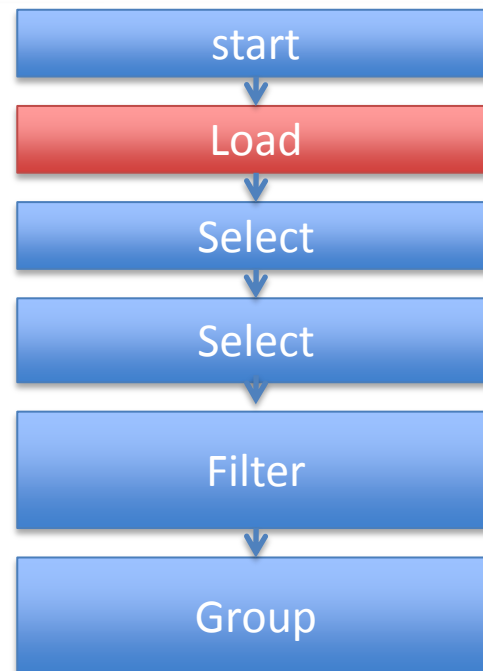
总结与问答

- 把用户编写的计算逻辑翻译为便于编译程序理解的中间码（语法树、数据流图）
- 前端代码运行一遍，产生结果是中间码
- 相当于编译技术中的parser
- 中间码用JSON表示

```
1 function get_site($rec){
2     $parts = explode( '/', $rec['_Url'] );
3     $rec['_Site'] = $parts[0];
4     return $rec;
5 }
6
7 $temp_view = DQuery::input()
8 ->select( array( '_Url','_Res' ) )
9 ->select( 'get_site' )
10 ->filter( array( '_Site', 'match' ) )
11
12 $temp_view
13 ->group( array( '_Site' ) )
14 ->each(
15     DQuery::count( '*', '_QueryCnt' )
16     DQuery::sum( '_Res', '_ResSum' )
17 )->outputAsFile( 'query_ad_shows',
18     'StorerUtils::jsonLine' );
```



```
[
  {
    "cmd": "load",
    "path": null
    "using": "SchemaReader"
    "from": 0
    "options": {"max_item_in_mem": 100000}
    "include": [25]
  }, {"cmd": "filter"...}, {"cmd": "group"...}, .....
```



```

query: SELECT fiel_list FROM form_list WHERE
clause
    fiel_list: fiel_list ' , ' fiel_name | fiel_name |
    ' * ' ;
    fiel_name: ID | FID ;
    form_list: form_list ' , ' form_name | form_name ;
    form_name: ID ;
    clause: clause AND clause | clause OR clause |
clause_express | ' ( ' clause ' ) ' ;
    clause_express: element compopera element ;
    element: element mathopera element | element1 |
    ' ( ' element ' ) ' ;
    compopera: ' < ' | ' > ' | ' = ' | ' > = ' | ' <
= ' ;
    mathopera: ' + ' | ' - ' | ' * ' | ' / ' ;
    element1: ID | ICON | FID | ' ' ' STR1 ' ' ' ;

```

## 语言定义（词法、语法分析）

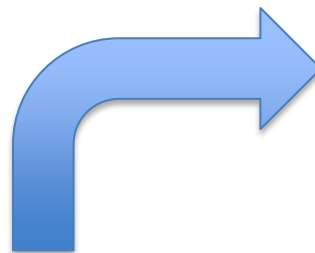


## 动作（生成中间码）

```

% %
SELECT | yylval. str = yytext; return (SELECT); |
";" | yylval. str = yytext; return (ENDMARK); |
/* ;为结束符 */
.....
| alpha | | alnum | * | yylval. str = yytext; return
(ID); | /* 变量名或字段各 */
| digit | + | yylval. str = yytext; return (I-
CON); | /* 常数 */
| alpha | | alnum | | . | | alpha | | alnum | * | yyl-
val. str = yytext; return (FID); |
/* 带表或视图名的字段名 */
| alnum | + | yylval. str = yytext; return (STR1); |
/* 字符串 */
[/\n]; /* 去除空格、制表符和换行 */
% %

```



```

[
{
    "cmd": "load",
    "path": null
    "using": "SchemaReader"
    "from": 0
    "options":
    {"max_item_in_mem":
    100000}
    "include": [25]
}
, {"cmd": "filter".....},
{"cmd": "join".....},..... .....
```

## 中间码



发展历程

一个例子

前端处理



中间语言翻译

运行时

总结与问答

## 正规化

- 将数据流图变成完整的方便后续处理的数据流图



## 算子替换

- 将实现复杂的算子等价替换成多个简单算子



## 优化

- 对数据流图进行各种优化，使执行效率提高



## 阶段划分(可选)

- 划分为多个MapReduce执行阶段
- 不划分会生成单机程序



## Schema推导、字段偏移量推导

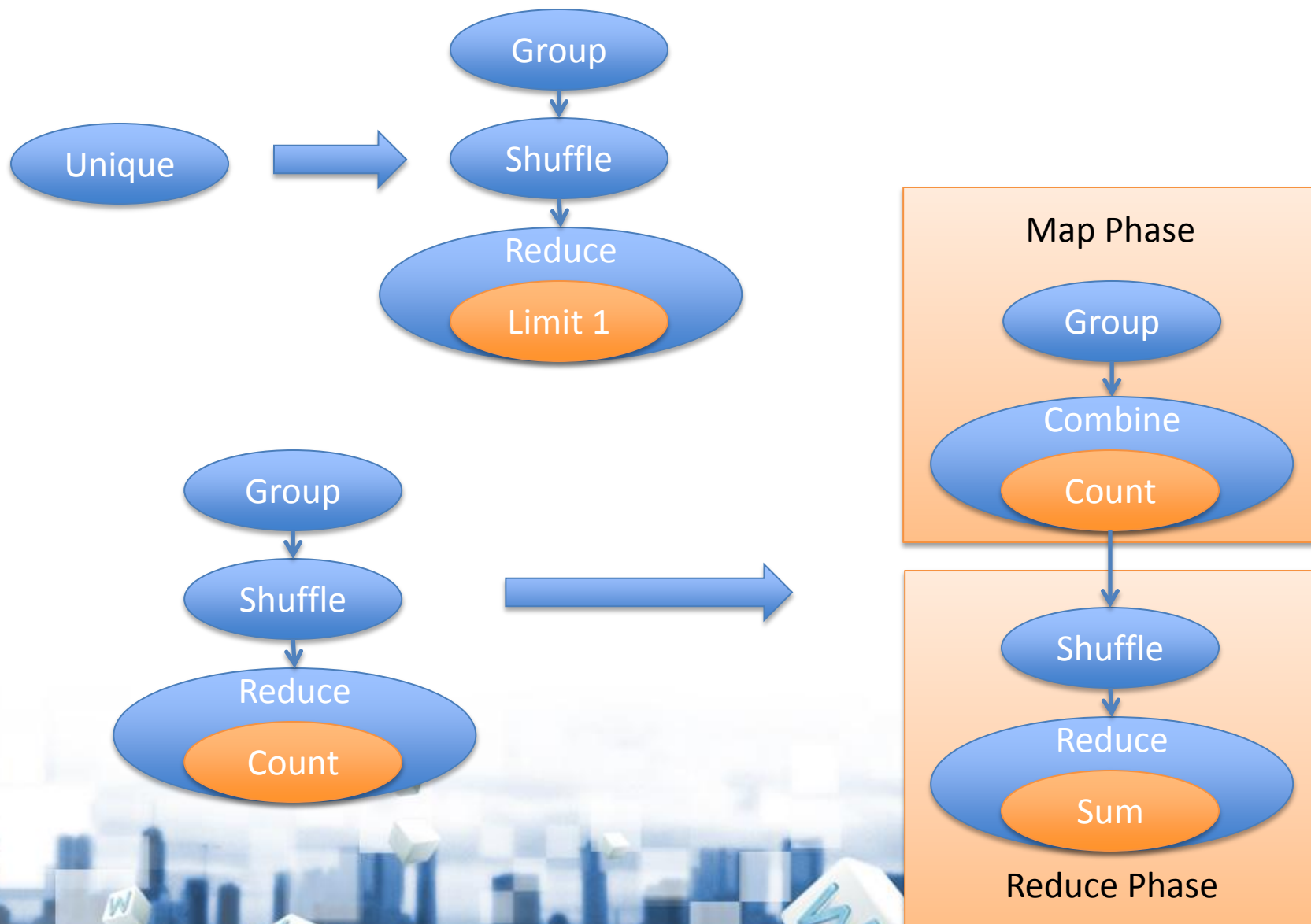
- 推导每一算子产出的表 schema，以及字段偏移量



## 代码生成

- 生成真正可执行的代码或数据流图





## ■ 多种优化策略

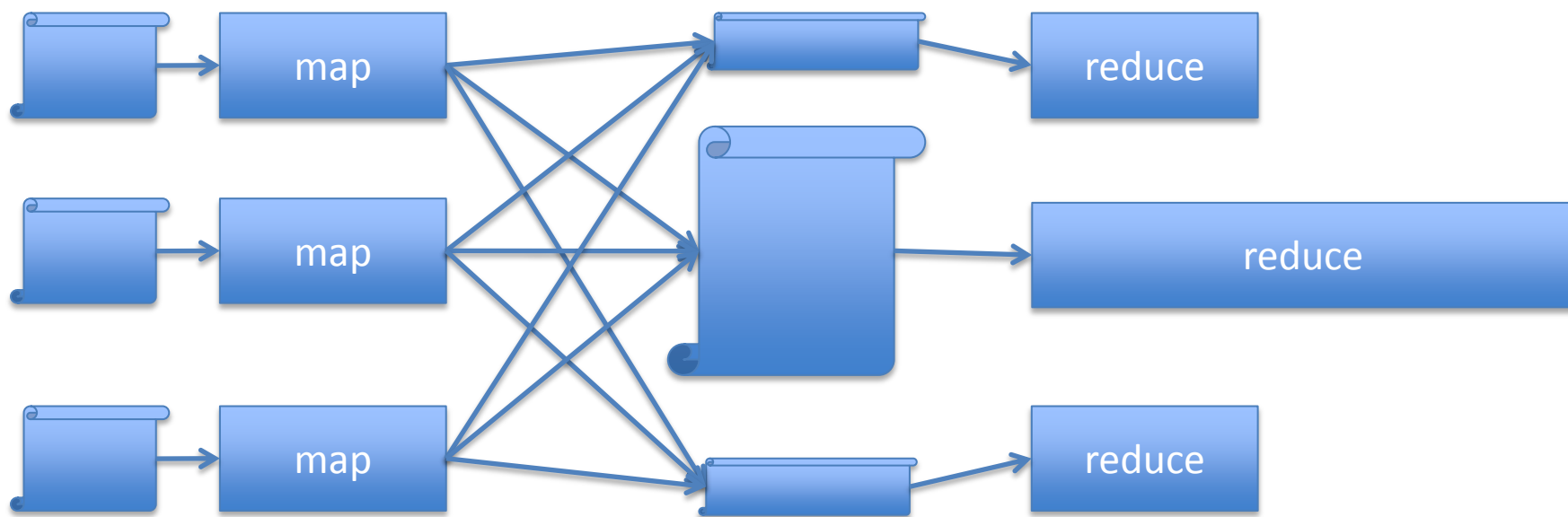
- 多任务合并
- 等价算子合并
- Combiner优化
- Cached Combiner优化
- 同key Join合并优化
- 公共子表达式提取
- .....

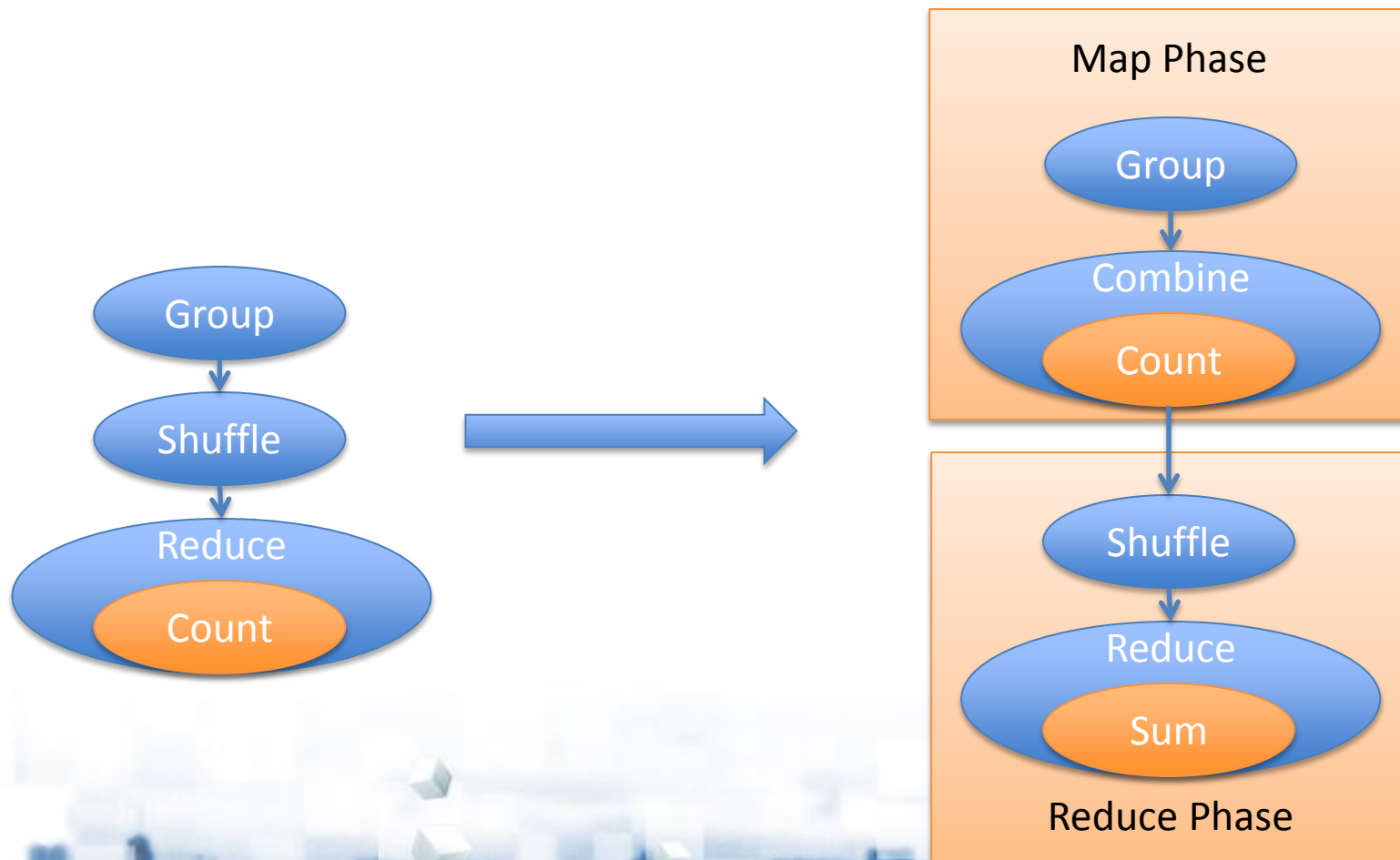
## ■ 核心思想

- 面向应用特点、减少作业轮数、减少I/O、减少重复计算

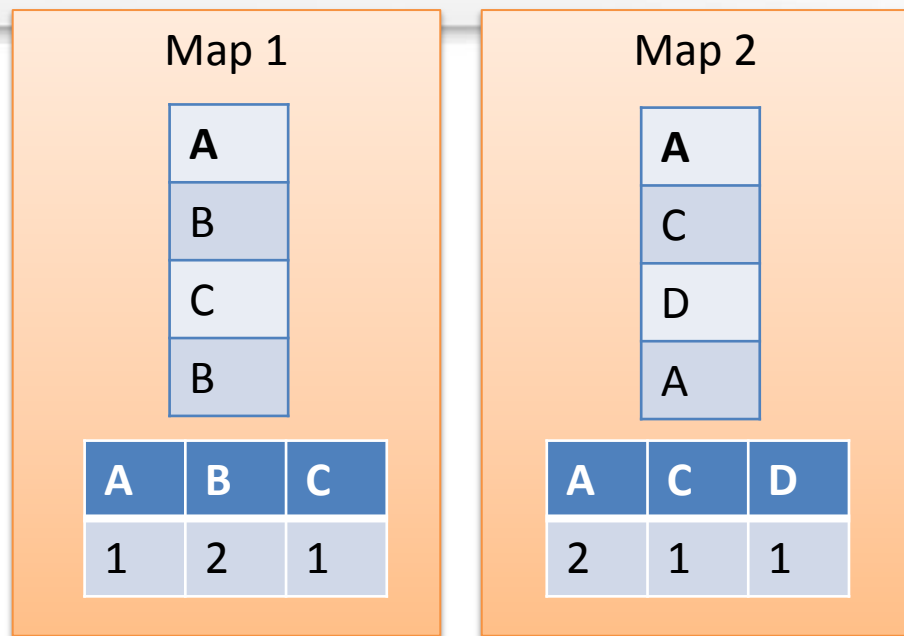
## 优化举例



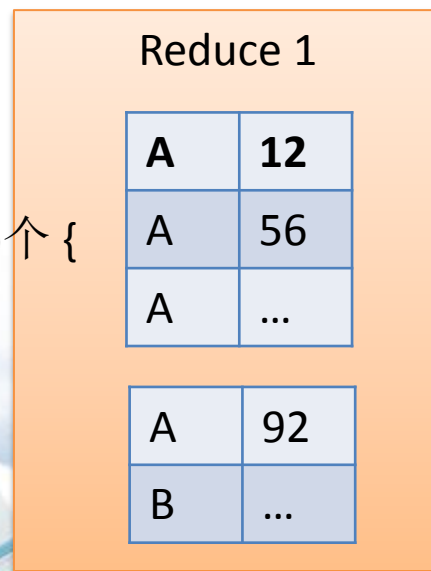




- 不采用系统Combiner
  - 额外磁盘I / O
- 编译到Mapper中实现
  - 内存中的hash字典



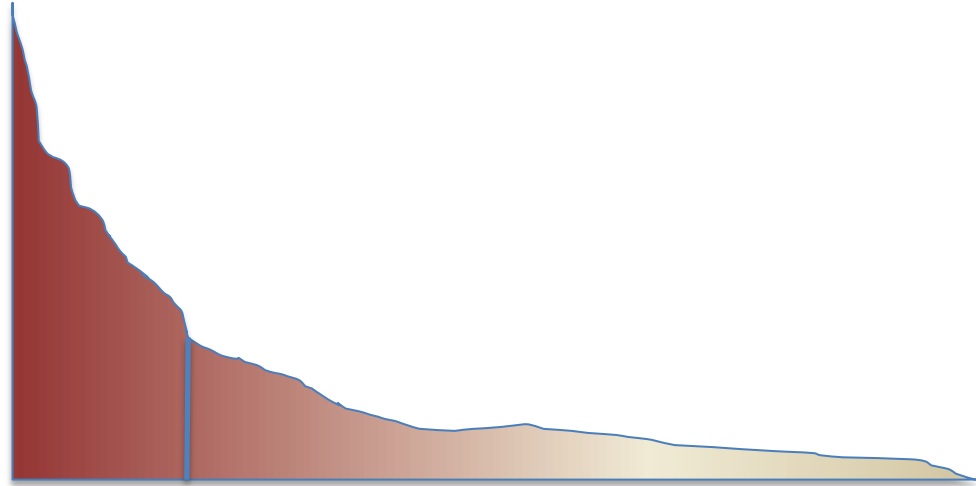
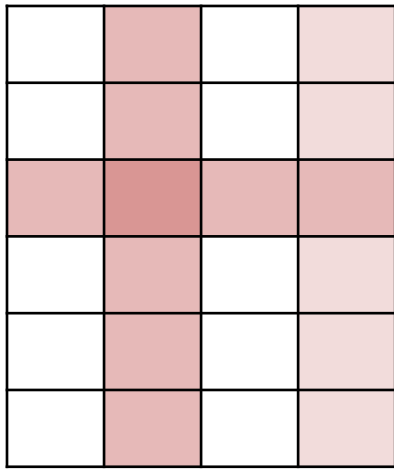
$O(\text{mapper}) \uparrow \{$



- |  |  |  |  |
|--|--|--|--|
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |

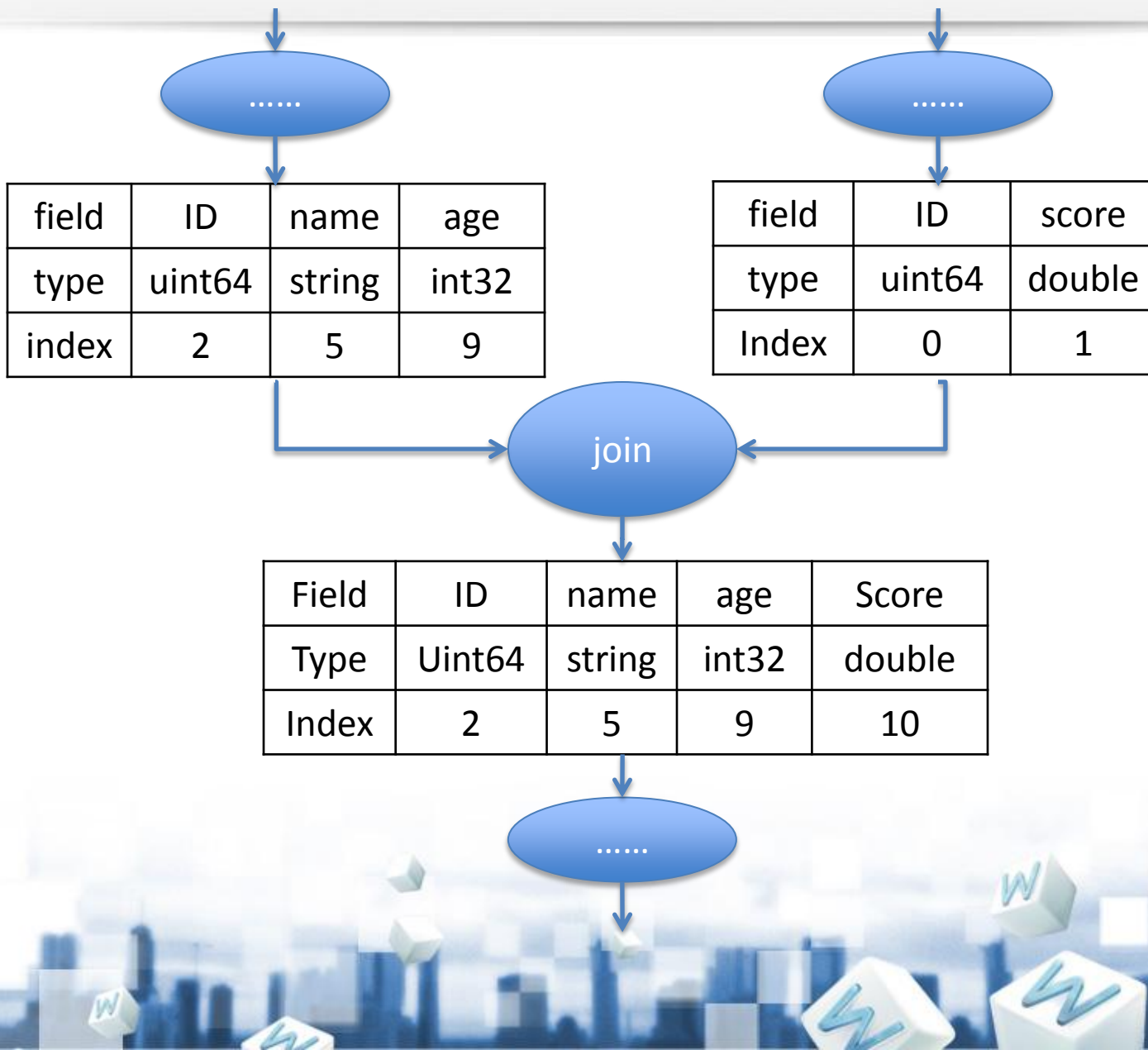


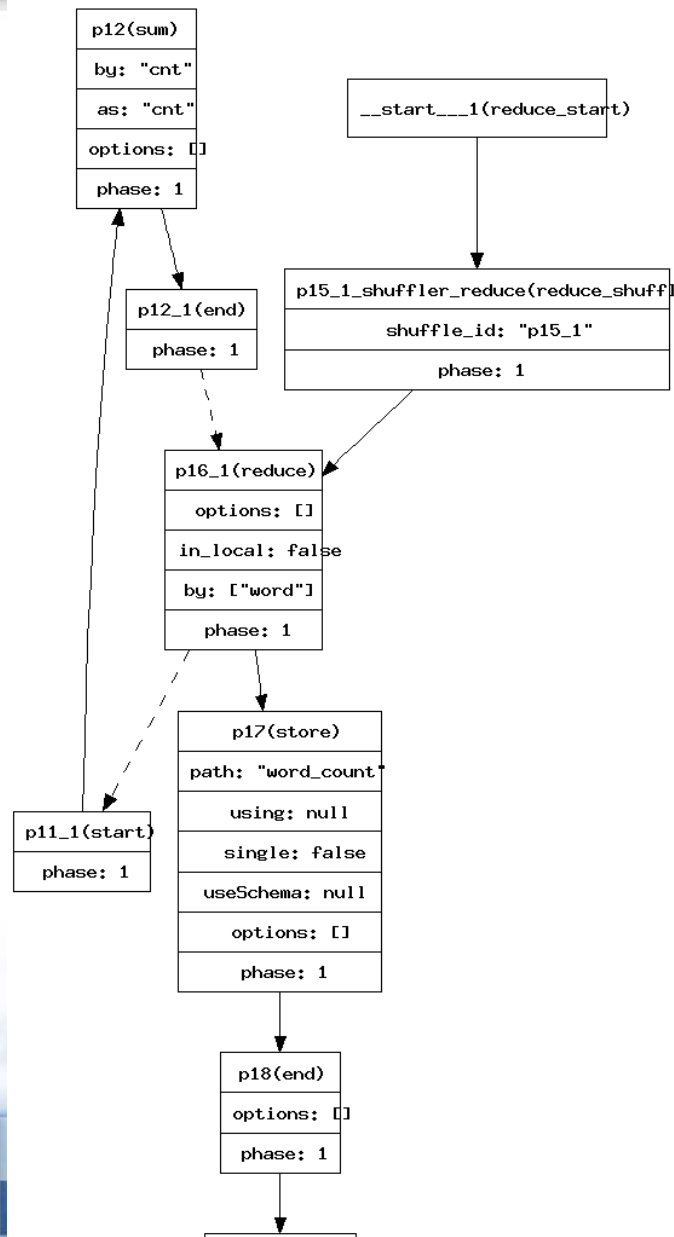
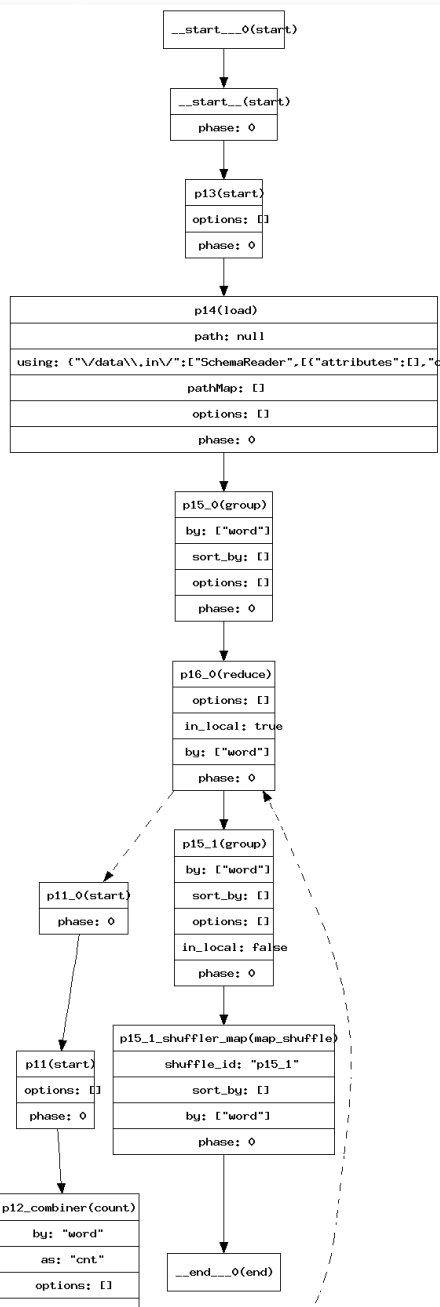




hot key	warm key	warm key	cool key	cool key	cool key







0.php (~/.doc/disql/tech\_salon/word\_count) - VIM

0.php (~/.doc/disql/tech\_salon/word\_count) - VIM 86

```

56 /***** vertex section *****/
57 $__start__ = new Starter;
58 $__start__->set_id('__start__');
59
60 $p13 = new Starter;
61 $p13->set_id('p13');
62
63 $p14 = new Loader(NULL,array (
64     '/data\\.in/' =>
65     array (
66         0 => 'SchemaReader',
67         1 =>
68         array (
69             0 =>
70             array (
71                 'attributes' =>
72                 array (
73                 ),
74                 'children' =>
75                 array (
76                     0 =>
77                     array (
78                         'attributes' =>
79                         array (
80                         ),
81                         'name' => 'word',
82                         'type' => 'string',
83                     ),
84                 ),
85                 'name' => 'res',
86                 'type' => 'struct',
87             ),
88         ),
89     ),
90 );
91 $p14->set_id('p14');
92
93 $p15_0 = new Grouper(array (
94     0 => 'word',
95 ));

```

```

111 $p12_combiner->set_to('p12');
112
113 $p12_0 = new Ender;
114 $p12_0->set_id('p12_0');
115
116 $p15_1 = new Grouper(array
117     0 => 'word',
118 );
119 $p15_1->set_id('p15_1');
120
121 $p15_1_shuffler_map = new
122     0 => 'word',
123 ),NULL,array (
124 );
125 $p15_1_shuffler_map->set_id
126
127 $__start__0 = new Starter;
128 $__start__0->set_id('__sta
129
130 $__end__0 = new Ender;
131 $__end__0->set_id('__end__
132
133 /***** edge section *****/
134 $__start__->follow($__start
135 $p13->follow($__start__);
136 $p14->follow($p13);
137 $p15_0->follow($p14);
138 $p16_0->follow($p15_0);
139
140 $p11->follow($p11_0);
141 $p12_combiner->follow($p11
142 $p12_0->follow($p12_combine
143 $p15_1->follow($p16_0);
144 $p15_1_shuffler_map->follow
145
146 $__end__0->follow($p15_1_s
147 /***** child section *****/
148 $p16_0->set_children($p11_0
149 /***** excution section ***
150 $__start__0->init(0);

```

0.cpp (~/.doc/disql/tech\_salon/word\_count) - VIM

0.cpp (~/.doc/disql/tech\_salon/word\_count) - VIM 78x41

```

1 #include "Starter__start__.h"
2 #include "Starter_p13.h"
3 #include "Loader_p14.h"
4 #include "Grouper_p15_1.h"
5 #include "MapShuffler_p15_1_shuffler_map"
6 #include "Starter__start__0.h"
7 #include "Ender__end__0.h"
8 #include <disql/types.h>
9 #include <disql/Context.h>
10 #include <disql/Dumper.h>
11 using namespace disql;
12
13 int main(int argc, char *argv[]){
14     Context::initialize(argc, argv);
15
16     bsl::Exception::set_line_delimiter("\n");
17     try{
18         /***** vertex section *****/
19         Starter__start__ * __start__ = new Starter__start__(1);
20         Starter_p13 * p13 = new Starter_p13(1);
21         Loader_p14 * p14 = new Loader_p14(2);
22         Grouper_p15_1 * p15_1 = new Grouper_p15_1(1);
23         MapShuffler_p15_1_shuffler_map * p15_1_shuffler_map = new MapShuffler_p15_1_shuffler_map(1);
24         Starter__start__0 * __start__0 = new Starter__start__0(1);
25         Ender__end__0 * __end__0 = new Ender__end__0(1);
26
27         /***** child section *****/
28
29         /***** edge section *****/
30         __start__->follow(*__start__0);
31         p13->follow(*__start__);
32         p14->follow(*p13);
33         p15_1->follow(*p14);
34         p15_1_shuffler_map->follow(*p15_1);
35
36         __end__0->follow(*p15_1_shuffler_map);
37
38         /***** edge section *****/
39         p12->follow(*p11_1);
40         p12_1->follow(*p12);
41         p16_1->follow(*p15_1_shuffler_map);
42         p17->follow(*p16_1);
43         p18->follow(*p17);
44         __end__->follow(*p18);
45         p15_1_shuffler_reduce->follow(*p15_1_shuffler_map);
46         __end__1->follow(*__end__);
47
48         /***** exec section *****/
49         bsl::var::IFactory& factory = C
50         Value main_record(factory);
51         __start__1->init(0);
52         __start__1->process(0, main_re
53         __start__1->fini(0);
54         /***** clean section *****/
55         delete p12;
56         delete p12_1;
57         delete p11_1;
58         delete p16_1;
59         delete p17;
60         delete p18;
61         delete __end__;
62         delete p15_1_shuffler_reduce;
63         delete __start__1;
64         delete __end__1;
65     } catch(bsl::Exception &e){
66         fprintf(stderr, "BSL EXCEPT
67     } catch(std::exception &e){
68         fprintf(stderr, "STD EXCEPT
69     } catch(...){
70         fprintf(stderr, "RUNTIME ER

```

发展历程

一个例子

前端处理

中间语言翻译



运行时

总结与问答

```
class  
Processor
```

```
init()
```

```
process()
```

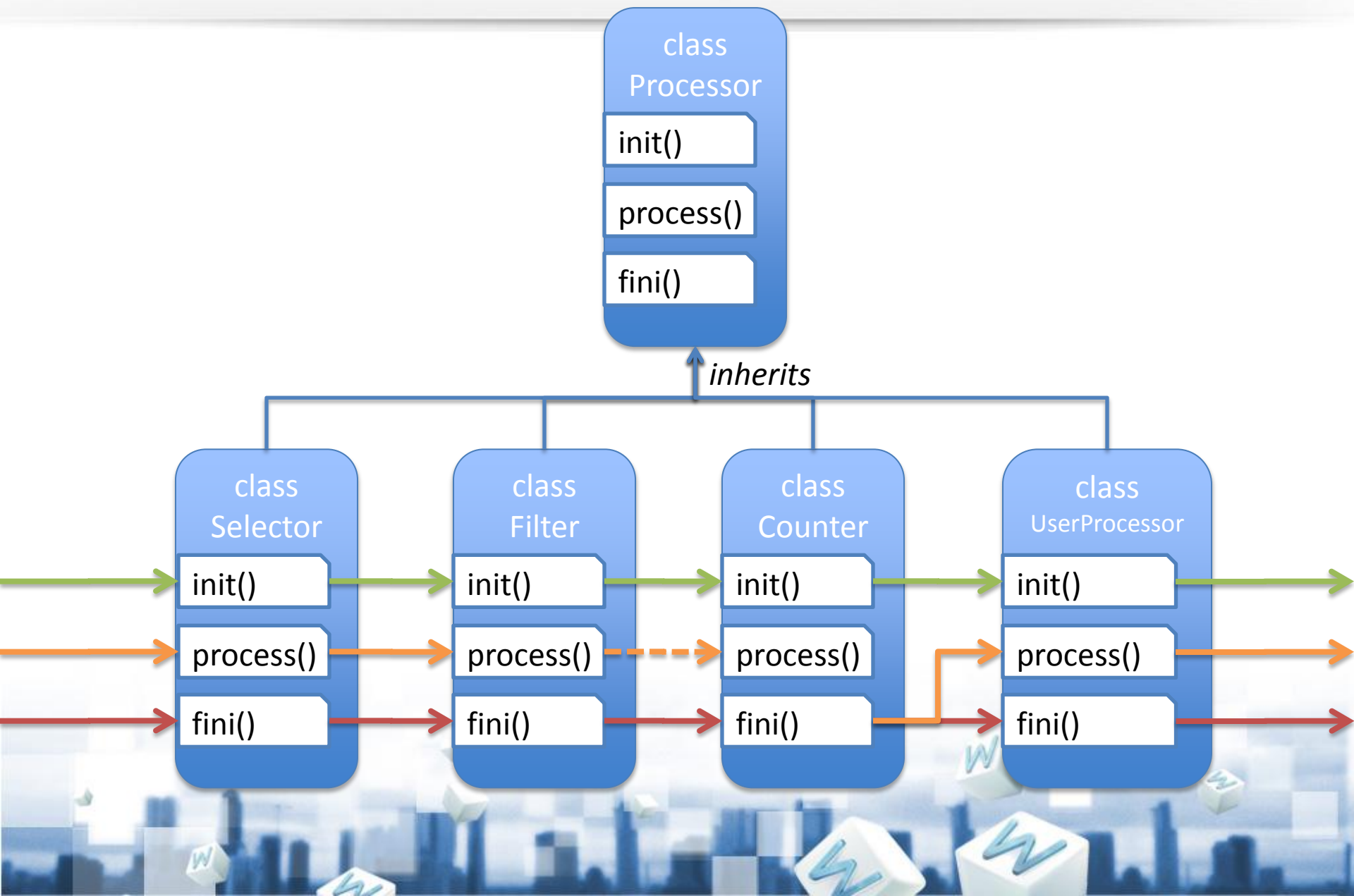
```
fini()
```

初始化一组数据处理

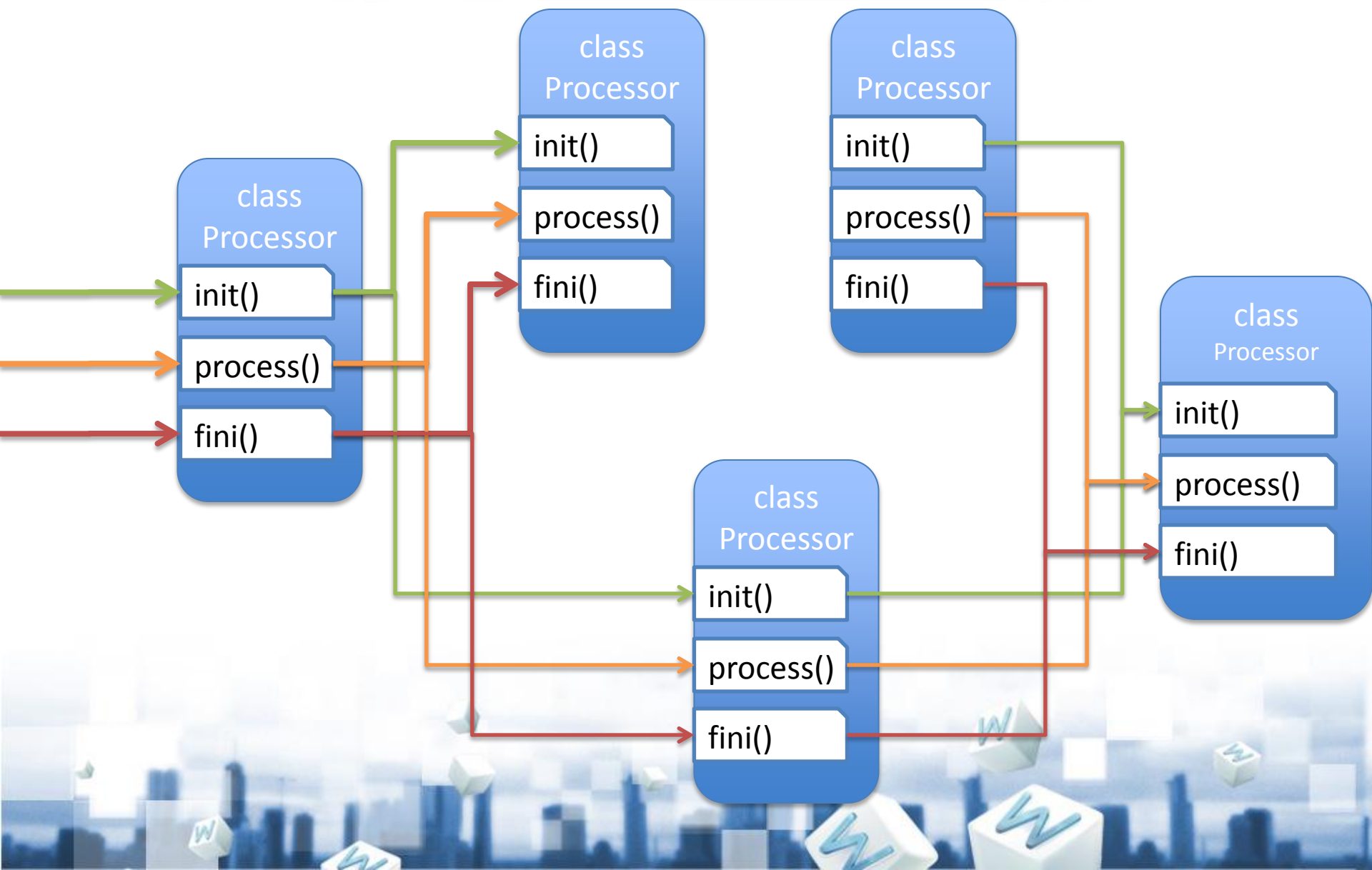
处理一组中的一条数据记录（多次调用）

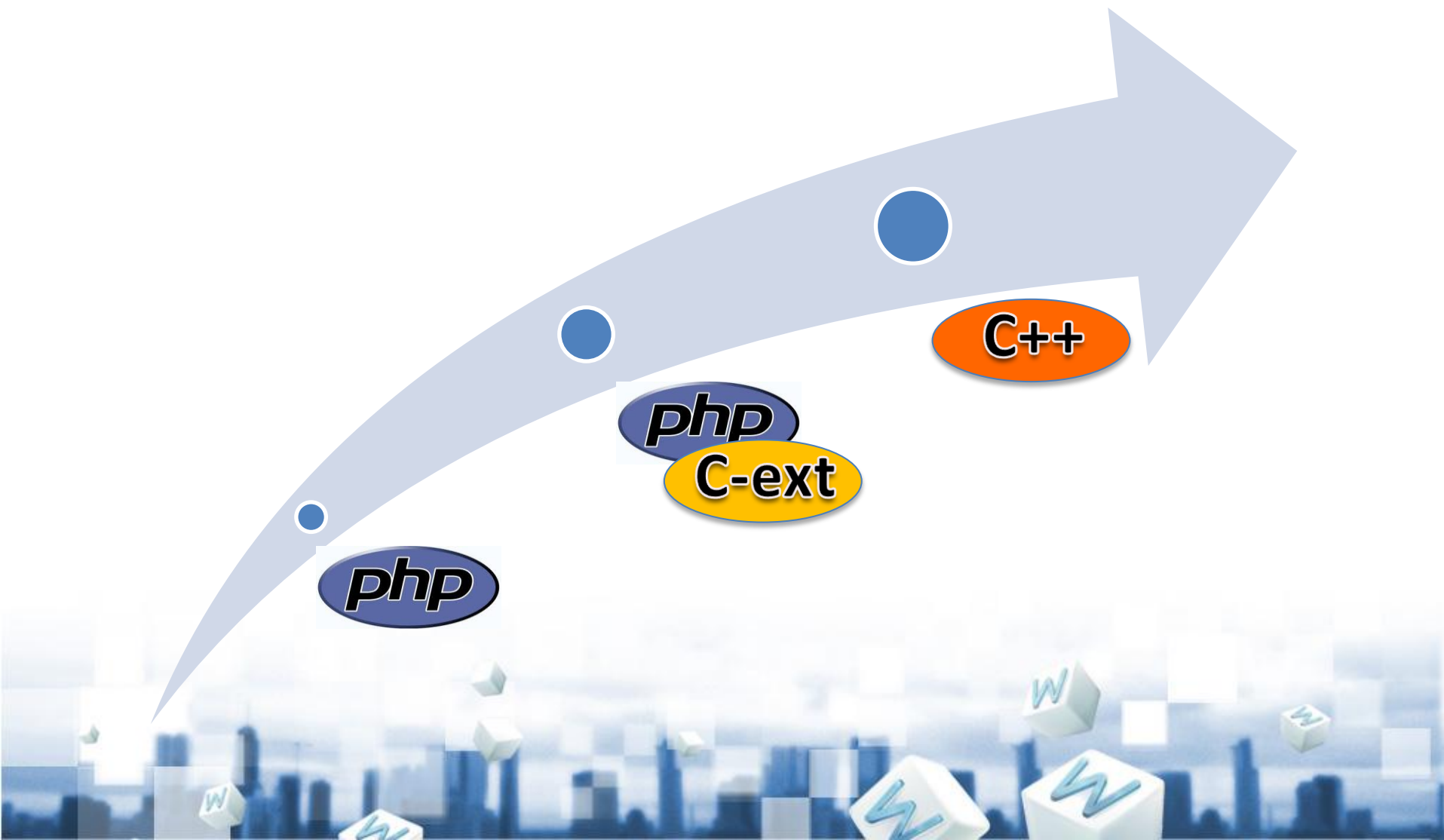
结束一组数据处理

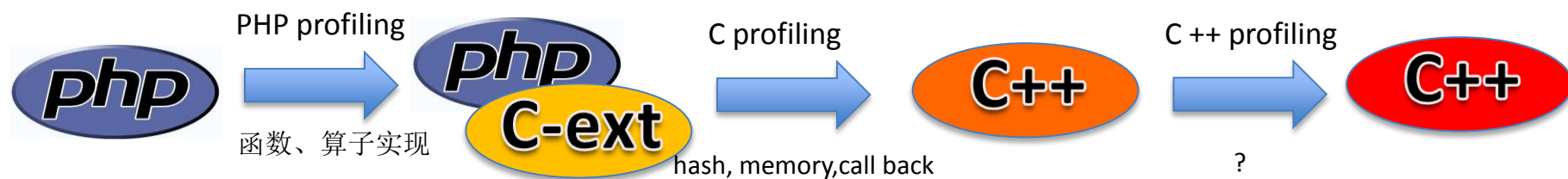
# 不同算子对应不同的Processor











- Hash操作
- Memory
- User Callback



PHP Runtime

```
record["user"]["name"]
```



C++ Runtime

下标推导

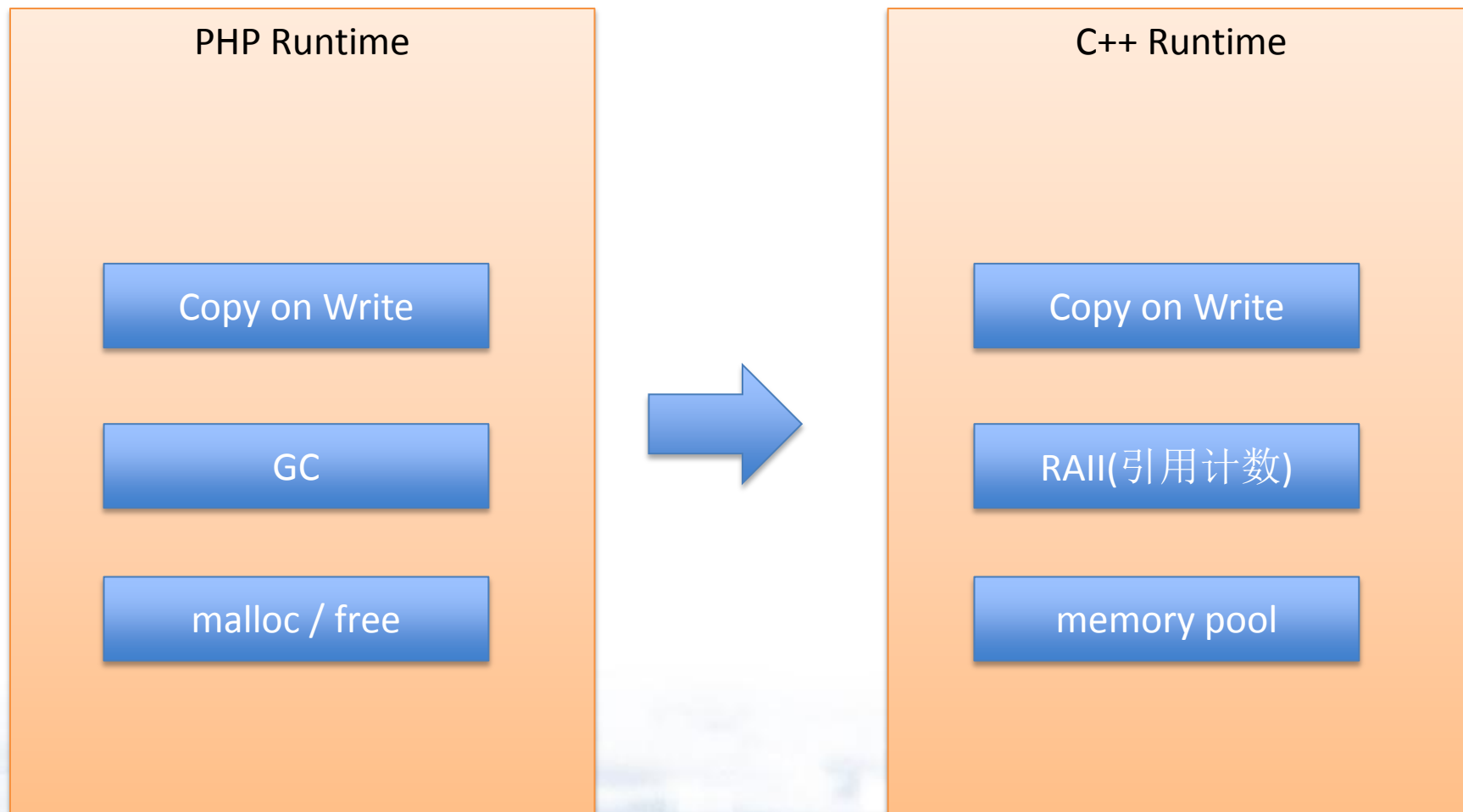
```
const int user= 1;  
const int name = 2;
```

用户代码、算子实现

```
record[user][name]
```

compile

```
record[1][2]
```



```
1 function expand($fields){
2     $ret=array();
3     for($i=0;$i<count($fields['_Dis'])[0];$i++) {
4         if($fields['_Dis'][0][$i]>0 && $fields['_Cn'][0][$i]!='baidu_fc_gusuan'
5         &&$fields['_Cn'][0][$i]!='baidufcidear_pg' && $fields['_Cn'][0][$i]!='baiduadrquery_pg')
6     {
7         for($j=0;$j<$fields['_Dis'][0][$i];$j++){
8             $r=array();
9             $r['_Srchid']=$fields['_S'];
10            $r['_Cmatch']=$fields['_Im_apres'][27][$i][$j];
11            $r['_Rank']=(string)$fields['_Absrk'][0][$i][$j];
12            $r['_query']=$fields['_Query'];
13            $r['_ip']=$fields['_Ip'];
14            $r['_cn']=$fields['_Cn'][0][$i];
15            $r['_bd']=$fields['_Bd'];
16            $r['_dis']=$fields['_Dis'][0][$i];
17            $r['_wd']=$fields['_Wd'][0][$i][$j];
18            $r['_pid']=$fields['_Pid'];
19            $r['_cid']=$fields['_Cid'];
20            $r['_tim']=$fields['_Tim'];
21            $r['_term']=$fields['_Term'][0][$i][$j];
22            $r['_pres_1']=$fields['_Pres_1'];
23            $r['_pn']=$fields['_Pn'];
24            $r['_eq']=$fields['_Eq'];
25            $r['_qs']=$fields['_Extra']['qs'];
26            $r['_const1']=1;
27            global $valfields;
28            foreach($valfields as $val)
29            {
30                if($r[$val]==' ' || $r[$val]==null)
31                    $r[$val]='-';
```



```
1 #include "disql_udf.h"
2 #include "disql_udf_schema.h"
3
4 extern "C" void expand(
5     bsl::var::LValue& fields, bsl::var::LValue& result, disql::CallbackContext& context) {
6     int k = 0;
7     if (!fields[ASP_LOG::_Dis].is_array()) return;
8     if (!fields[ASP_LOG::_Dis][0].is_array()) return;
9     int dis0_size = fields[ASP_LOG::_Dis][0].size();
10    for (int i = 0; i < dis0_size; i++) {
11        int dis0i;
12        const char * cn;
13        if (fields[ASP_LOG::_Dis][0][i].is_null() || strcmp(fields[ASP_LOG::_Dis][0]
14[i].c_str(), "") == 0) dis0i = 0;
15        else dis0i = fields[ASP_LOG::_Dis][0][i].to_int32();
16        if (fields[ASP_LOG::_Cn][0][i].is_null()) cn = "";
17        else cn = fields[ASP_LOG::_Cn][0][i].c_str();
18        if (dis0i > 0 && strcmp("baidu_fc_gusuan", cn) != 0 && strcmp("baidufcidear_pg", cn)
19!= 0 && strcmp("baiduadrquery_pg", cn) != 0) {
20            for (int j = 0; j < dis0i; j++) {
21                result[k][EXPAND::_Srchid] = fields[ASP_LOG::_S];
22                result[k][EXPAND::_Cmatch] = fields[ASP_LOG::_Im_apres][27][i][j];
23                result[k][EXPAND::_Rank] = fields[ASP_LOG::_Absrk][0][i][j];
24                result[k][EXPAND::_query] = fields[ASP_LOG::_Query];
25                result[k][EXPAND::_ip] = fields[ASP_LOG::_Ip];
26                result[k][EXPAND::_cn] = fields[ASP_LOG::_Cn][0][i];
27                result[k][EXPAND::_bd] = fields[ASP_LOG::_Bd];
28                result[k][EXPAND::_dis] = fields[ASP_LOG::_Dis][0][i];
29                result[k][EXPAND::_wd] = fields[ASP_LOG::_Wd][0][i][j];
30                result[k][EXPAND::_pid] = fields[ASP_LOG::_Pid];
31                result[k][EXPAND::_cid] = fields[ASP_LOG::_Cid];
```

- 依赖百度在C++编程上的丰富积累，尤其是BSL库
  - BSL = Baidu Standard Library，使用最广泛的基础库
  - 各种高效的泛型容器
  - 内存池
  - 不弱于Java异常的异常支持
  - 动态类型支持
    - 像写PHP那样编写C++代码
  - 格式转换
    - JSON、内部二进制格式.....
  - 语言扩展支持
    - PHP、Python、Perl.....



发展历程

一个例子

前端处理

中间语言翻译

运行时

总结与问答



- 前端
  - DQuery: 2037
- 中间语言翻译
  - 4552
- 运行时
  - PHP Runtime: 7283
    - PHP: 2842
    - C Extension: 4441
  - C++ Runtime: 8318



- 轻量级
- 类SQL逻辑(非常简约)
  - 封装所有SQL算子的M/R分布式实现：
    - 分组、聚合、表连接、行列过滤、集合操作、输入输出格式转换
- DAG数据流
  - 自动翻译为一轮或多轮MapReduce
  - 也可翻译为单机计算或数据流图
- 逻辑顺序而非SQL顺序
- 支持PHP自定义函数（简洁）
- 支持C++自定义函数（同样简洁+高效）和C-Runtime **NEW!**
  - 全自动高效内存管理（RAII + 内存池）
  - 廉价对象复制（Copy On Write）
  - 字段操作翻译为数组操作，无字典查找（schema推导）
  - C++的性能，PHP的开发代价！

## ■ 分析程序数量增长

程序类型	4月1日	10月27日	增长	增长百分比
简单编辑	3540	4761	1221	+34.5%
DQuery模式	1153	3359	2206	<b>+191%</b>
复杂编辑	1569	2963	1394	+88.9%

## ■ 分析程序输入量

程序类型	占比
Web表单	24%
DQuery模式	43%
裸MapReduce	33%

} 67%

## ■ 用户数量

角色	数量	占比
PM	1352	47.4%
RD	1174	41.2%
OP	190	6.66%
其他	136	4.77%
总数	2852	100%

# Q & A

更多问题，可以向微博[@陈晓鸣在百度](#)提问





谢谢！

