



## 百度日志分析技术分享



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)

# 互联网公司技术架构系列资料



为您悉心整理

/\* 让工作重新关于成长和成就、关于快乐和分享、关于梦想和荣光 \*/



LOG中自有黄金屋

日志分析基本过程

百度日志分析成长历程

深入LSP平台

深入DISQL语言

总结与问答

■ 46.70.93.94 - - [11/Nov/2011:11:11:11 -1100] "GET /book/1984.html HTTP/1.1 "404 2326  
http://www.baidu.com/s?wd=1984&rsv\_bp=0&rsv\_spt=3  
&inputT=947 "Mozilla/5.0(iPad; U; CPU iPhone OS 3\_2  
like Mac OS X; en-us) AppleWebKit/531.21.10 (KHTML,  
like Gecko) Version/4.0.4 Mobile/7B314 Safari/531.21.10  
"



- 46.70.93.94 - -
- [11/Nov/2011:11:11:11 -1100]
- "GET /book/1984.html HTTP/1.1"
- 404
- 2326
- "http://www.baidu.com/s?wd=1984&rsv\_bp=0&rsv\_spt=3&inputT=947"
- "Mozilla/5.0(iPad; U; CPU iPhone OS 3\_2 like Mac OS X; en-us) AppleWebKit/531.21.10 (KHTML, like Gecko) Version/4.0.4 Mobile/7B314 Safari/531.21.10 "



- 46.70.93.94 - -
- [11/Nov/2011:11:11:11 -1100]
- "GET /book/1984.html HTTP/1.1"
- 404
- 2326
- "http://www.baidu.com/s?wd=1984&rsv\_bp=0&rsv\_spt=3&inputT=947"
- "Mozilla/5.0(iPad; U; CPU iPhone OS 3\_2 like Mac OS X; en-us) AppleWebKit/531.21.10 (KHTML, like Gecko) Version/4.0.4 Mobile/7B314 Safari/531.21.10 "





- 46.70.93.94 - -
- [11/Nov/2011:11:11:11 -1100]
- "GET /book/1984.html HTTP/1.1"
- 404
- 2326
- "http://www.baidu.com/s?wd=1984&rsv\_bp=0&rsv\_spt=3&inputT=947"
- " Mozilla/5.0(iPad; U; CPU iPhone OS 3\_2 like Mac OS X; en-us) AppleWebKit/531.21.10 (KHTML, like Gecko) Version/4.0.4 Mobile/7B314 Safari/531.21.10"



- 46.70.93.94 - -
- [11/Nov/2011:11:11:11 -1100]
- GET /book/1984.html HTTP/1.1
- 404
- 2326
- "http://www.baidu.com/s?wd=1984&rs  
v\_bp=0&rsv\_spt=3&inputT=947 "
- "Mozilla/5.0(iPad; U; CPU iPhone OS  
3\_2 like Mac OS X; en-us)  
AppleWebKit/531.21.10 (KHTML,  
like Gecko) Version/4.0.4  
Mobile/7B314 Safari/531.21.10 "

Traffic to etsy.com

(Unique Visitors and Page Views in Millions - Oct '08 to Oct '09)

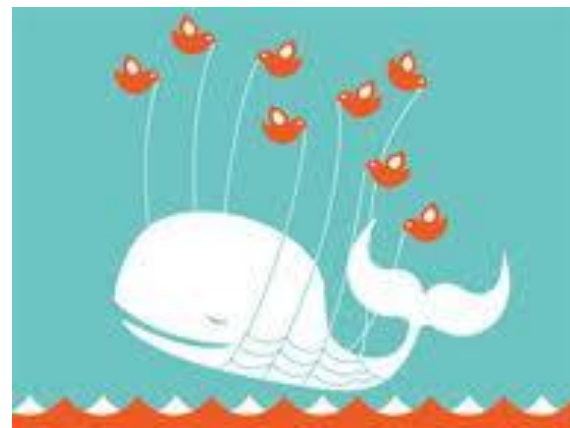




- 46.70.93.94 - -
- [11/Nov/2011:11:11:11 -1100]
- "GET /book/1984.html HTTP/1.1"
- 404
- 2326
- "http://www.baidu.com/s?wd=1984&rsv\_bp=0&rsv\_spt=3&inputT=947"
- " Mozilla/5.0(iPad; U; CPU iPhone OS 3\_2 like Mac OS X; en-us) AppleWebKit/531.21.10 (KHTML, like Gecko) Version/4.0.4 Mobile/7B314 Safari/531.21.10 "



- 46.70.93.94 - -
- [11/Nov/2011:11:11:11 -1100]
- "GET /book/1984.html HTTP/1.1"
- **404**
- 2326
- " http://www.baidu.com/s?wd=1984&rsv\_bp=0&rsv\_spt=3&inputT=947 "
- "Mozilla/5.0(iPad; U; CPU iPhone OS 3\_2 like Mac OS X; en-us) AppleWebKit/531.21.10 (KHTML, like Gecko) Version/4.0.4 Mobile/7B314 Safari/531.21.10 "









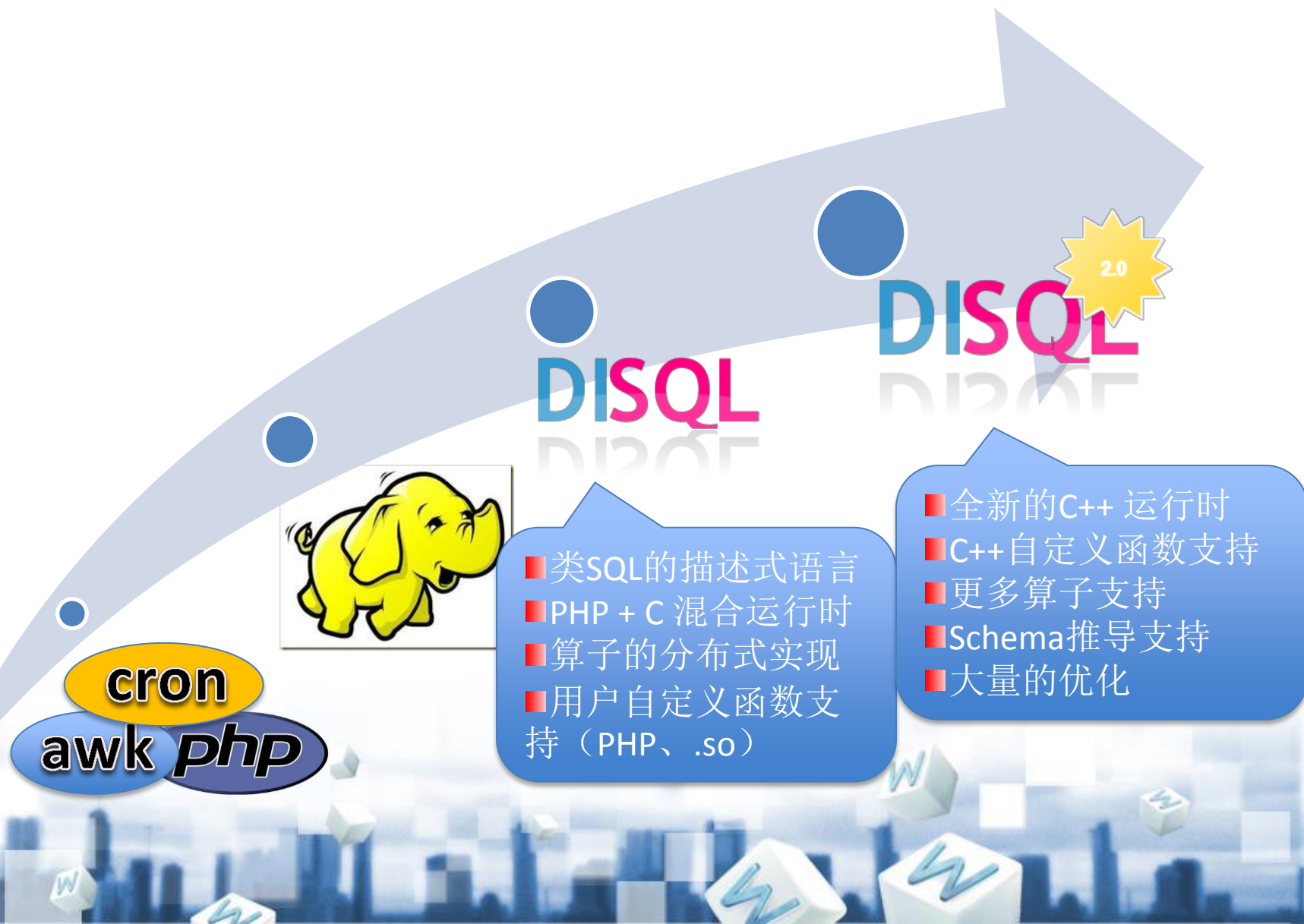
- LOG中自有黄金屋
- ➔ 日志分析基本过程
- 百度日志分析成长历程
- 深入LSP平台
- 深入DISQL语言
- 总结与问答





- LOG中自有黄金屋
- 日志分析基本过程
- ➡ 百度日志分析成长历程
- 深入LSP平台
- 深入DISQL语言
- 总结与问答





LOG中自有黄金屋

日志分析基本过程

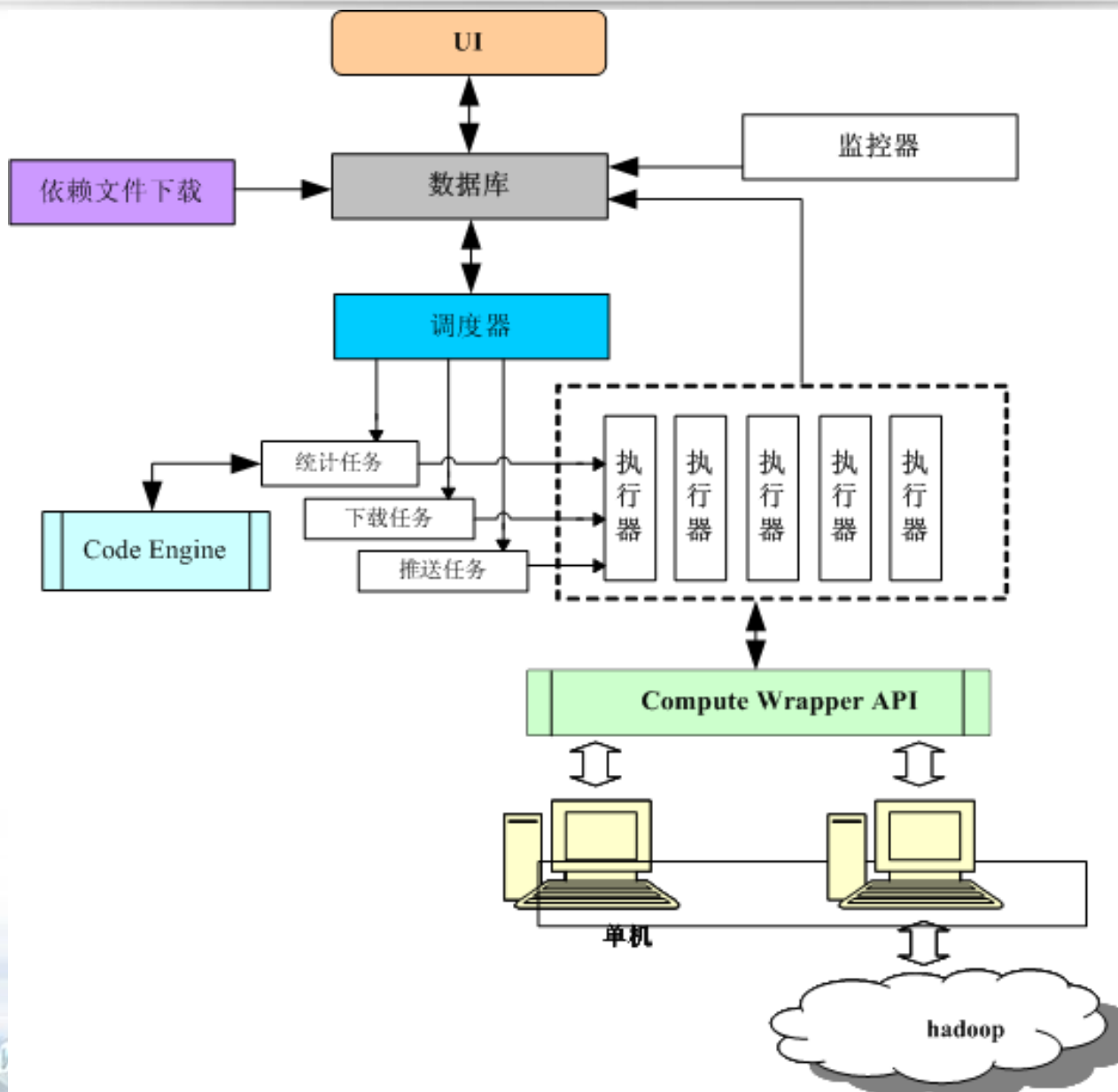
百度日志分析成长历程



深入LSP平台

深入DISQL语言

总结与问答



退出chenxiaomi

报表管理

我的收藏的报表

我创建的报表

所有报表

所有统计项

创建报表

统计管理

查看统计列表 | 新增

数据块管理

依赖文件管理 | 新增

查看项目列表 | 新增

任务管理

统计任务 | 今日

推送任务 | 今日

下载任务 | 今日

缓存任务 | 今日

投影任务 | 今日

启动统计任务

数据管理

系统帮助

产品线管理

运行状态监测

监控报警

版本: V1.0 [中文 | 日本語]

你当前的位置: [统计管理] - [统计列表]

运行选中的统计 运行搜索出来的所有统计 函数列表 函数新增

<input type="checkbox"/>	序号	产品线	创建者	统计名称	统计机房	类别	优先级	当前状态	基本操作
<input type="checkbox"/>	14454	bae	zhangwei07	hao123][二级]桌面快捷方式统计	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14453	wise	li_kai	携带手机号的PV数量	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14452	ns	yangran	[Openapp]平台首页点击分布	全部	复杂编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14451	iknow	tianchao	[活动任务_高考倒计时]效果_taskd	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14450	iknow	tianchao	[活动任务_高考倒计时]效果统计_taskui	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14449	ecom_holmesliangjiangzhang		tmp_site_keyword_count	全部	DQuery模式	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14448	skycn	cuisongyan	dddd	全部	简单编辑	3	未填写统计方法	编辑 复制 删除 更多>>
<input type="checkbox"/>	14447	bae	zhangwei07	[hao123][二级]tongzhi统计	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14446	wise	liujingwei02	wisebusiness_phonerecall_survey	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14445	ecom_cpro	suerfeng	fc_site_adtrade_stat_hourly_new	全部	DQuery模式	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14444	ubs	liushaowei	userdata_cookie_image	全部	复杂编辑	7	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14443	ecom_im	lvqicong	imdp_hourly_join_second	全部	DQuery模式	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14442	ecom_cpro	liuxudong	liuxudong_cmetstat_0_eclick_check	全部	DQuery模式	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14441	ubs	liushaowei	userdata_cookie_map_total	全部	复杂编辑	7	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14440	ubs	yuqi	newcookiesortqt	全部	复杂编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14439	exp	zuchang	[经验分享]每日原创经验提交条数	全部	简单编辑	3	未填写统计方法	编辑 复制 删除 更多>>
<input type="checkbox"/>	14438	ubs	linan03	userdata_ps_queryTitleLog	全部	简单编辑	3	未填写统计方法	编辑 复制 删除 更多>>
<input type="checkbox"/>	14437	space	luweichao	长连接前端统计	全部	DQuery模式	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14436	space	yangzhanye	aaa_[百度空间]认证用户统计	全部	简单编辑	3	有效	编辑 复制 删除 更多>>
<input type="checkbox"/>	14435	ubs	qumingcheng	temp_userdata_random_ip_query	全部	复杂编辑	6	有效	编辑 复制 删除 更多>>

共有 8827 条记录, 当前第 1/442 页

首页 上一页 下一页 尾页 转到第 页 转

统计查询

所属产品线:

不限

名称:

区分机房:

不限

类别:

不限

优先级:

不限

运行周期:

不限

创建者:

修改者:

创建时间:



# 编辑模式



[退出]chenxiaoming

## ■ 报表管理

■ 我收藏的报表

■ 我创建的报表

■ 所有报表

■ 所有统计项

■ 创建报表

## ■ 统计管理

■ 查看统计列表 | 新增

■ 日志模块管理 | 新增

■ 依赖文件管理 | 新增

■ 查询下载记录

## ■ 任务管理

■ 统计任务 | 今日

■ 下载任务 | 今日

■ 缓存任务 | 今日

■ 推送任务 | 今日

■ 启动统计任务

■ 调度器管理

## ■ 数据管理 ?

■ 数据推送 | 新增

■ 推送记录查询

■ 统计项保存时间设置

■ 关联关系列表 | 新增

## ■ 运行状态监测

■ 统计任务监控

■ 报表监控

您当前的位置: ecom\_cpro - test

切换编辑模式 ▾

检查

测试

提交

☒ 日志行实例 ☐ 上传文件测试[+ 添加数据源](#) [显示高级选项](#)1. 日志: ecom\_pb\_log [删除](#)

## 统计编辑区

[+ 计数](#)[+ Top统计](#)[+ 去重计数](#)[+ 数据过滤](#)[+ 更多](#)[全部收起](#)[帮助](#)

## 1. 分地域计数

\*分城市计数变量名  \*中文名 \*分省份计数变量名  \*中文名 \*IP字段: 

请选择代表"IP地址"的字

[+ 添加限制条件](#) [使用公共限制条件组](#)1. 

正则匹配

2. 

数字大于

为该统计添加注释(选填): [加和](#)  
[去重列表](#)  
[分时段计数](#)  
[分时段加和](#)  
[分时段去重](#)  
[分组计数](#)  
[分地域计数](#)  
[分地域去重](#)  
[新增记录统计](#)  
[重复记录统计](#)  
[Transmit统计](#)

## 公共限制条件组

[+ 添加公共限制条件组](#) [全部收起](#)1. [删除](#)[+ 添加限制条件](#)

我要测试

[测试返回信息](#) [取消滚动](#)

:) 您还没进行测试!

您当前的位置: cpro - doris\_log\_dcharge\_sf生成字典

切换编辑模式 ▾

统计数据源列表

[+ 添加数据源](#) [显示高级选项](#)类型: 日志 ▾ 数据源: doris\_log\_dcharge\_sf ▾ 日期: 当日 ▾ [删除](#)

字段列表

原始记录集的所有字段列表: \_Balance ▾ \_Balance

代码编辑区

[API手册](#) [帮助](#) [查看流程图](#)

```
1 DQuery::input ()
2 ->filter( array( array( '_Cmatch', '===', 202 ) ) )
3 ->select( array( '_Srcid', '_Descid', '_Bid' ) )
4 ->group( array( '_Srcid', '_Descid' ) )
5 ->each(
6     DQuery::count('*', '_Click' ),
7     DQuery::sum('_Bid', '_Gain' )
8 )
9 ->outputAsFile('doris_log_dict','doris日志生成字典',
10 'StorerUtils::phpLine',true);
```

检查

测试

提交

☒ 日志行实例 ☐ 上传文件测试

```
101      1947467 2055624
11184222
269167736      周星驰
206344 28176565
0.68      0.54
32918.62
^ 00005510      00000
```

我要测试

测试返回信息

**doris\_log\_dict:**

a:5:{s:7:"\_Srcid";s:16:"5bef9a5e"



LOG中自有黄金屋

日志分析基本过程

百度日志分析成长历程

深入LSP平台



深入DISQL语言

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



```
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', '/news\[^\.\]+\\.cn/' ) );
11
12 $temp_view
13 ->group( array( '_Site' ) )
14 ->each(
15     DQuery::count( '*', '_QueryCnt' ),
16     DQuery::sum( '_Res', '_ResSum' )
17 ) ->outputAsFile( 'query_ad_shows', '分站点访问量与广告量统计', 'StorerUtils::jsonLine' );
```

■ 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

■ 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

```
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];
```

- 用PHP表达的类SQL逻辑(非常简约)
  - 封装所有SQL算子的M/R分布式实现：
    - 分组、表连接、行列过滤、集合操作、输入输出格式转换
- 通过连续函数调用表达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
- 用逻辑顺序而非SQL顺序表达逻辑
- 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！

前端语  
言处理

中间语  
言翻译

运行时

- 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

```
[  
  {  
    "cmd": "load",  
    "path": null  
    "using": "SchemaReader"  
    "from": 17  
    "options": {"max_item_in_mem": 100000}  
    "include": [25]  
  },  
  {"cmd": "filter".....}, {"cmd": "join".....}, .....  
]
```



```

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); | /* 变量名或字段名 */
| 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 */
%%

```

```

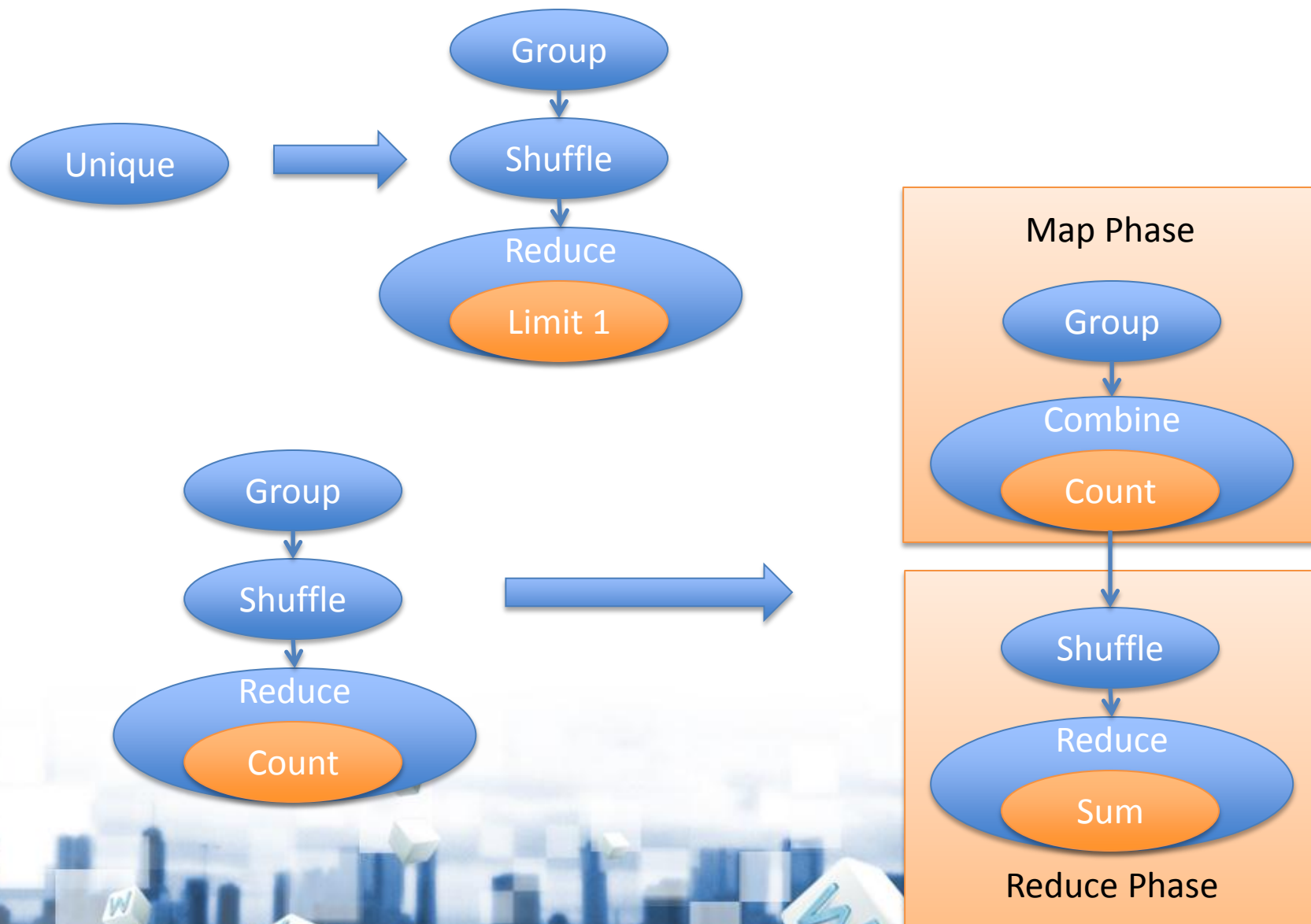
[
{
    "cmd": "load",
    "path": null
    "using": "SchemaReader"
    "from": 17
    "options":
    {"max_item_in_mem":
    100000}
    "include": [25]
}
, {"cmd": "filter".....},
{"cmd": "join".....},..... ..
]

```

## 中间码

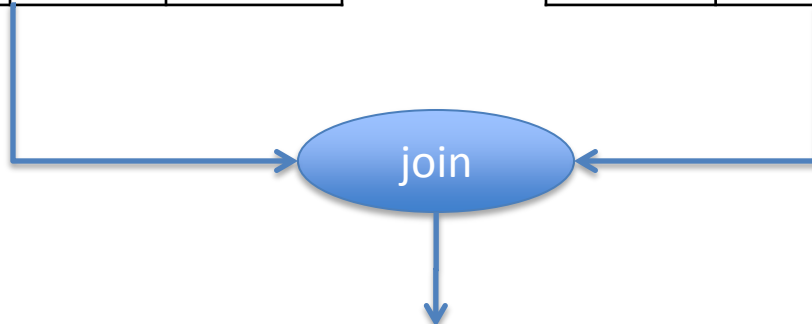


- 对数据流图作多次等价变换
  - 正规化
    - 将数据流图变成完整的方便后续处理的数据流图
  - 算子替换
    - 将实现复杂的算子等价替换成多个简单算子
  - 优化
    - 对数据流图进行各种优化，使执行效率提高
  - 阶段划分(可选)
    - 划分为多个MapReduce执行阶段
  - Schema推导、字段偏移量推导
    - 推导每一算子产出的表schema，以及字段偏移量
  - 代码生成 ( C++、PHP、DOT )
    - 生成真正可执行的代码



field	ID	name	age
type	uint64	string	int32
index	2	5	9

field	ID	score
type	uint64	double
Index	0	1

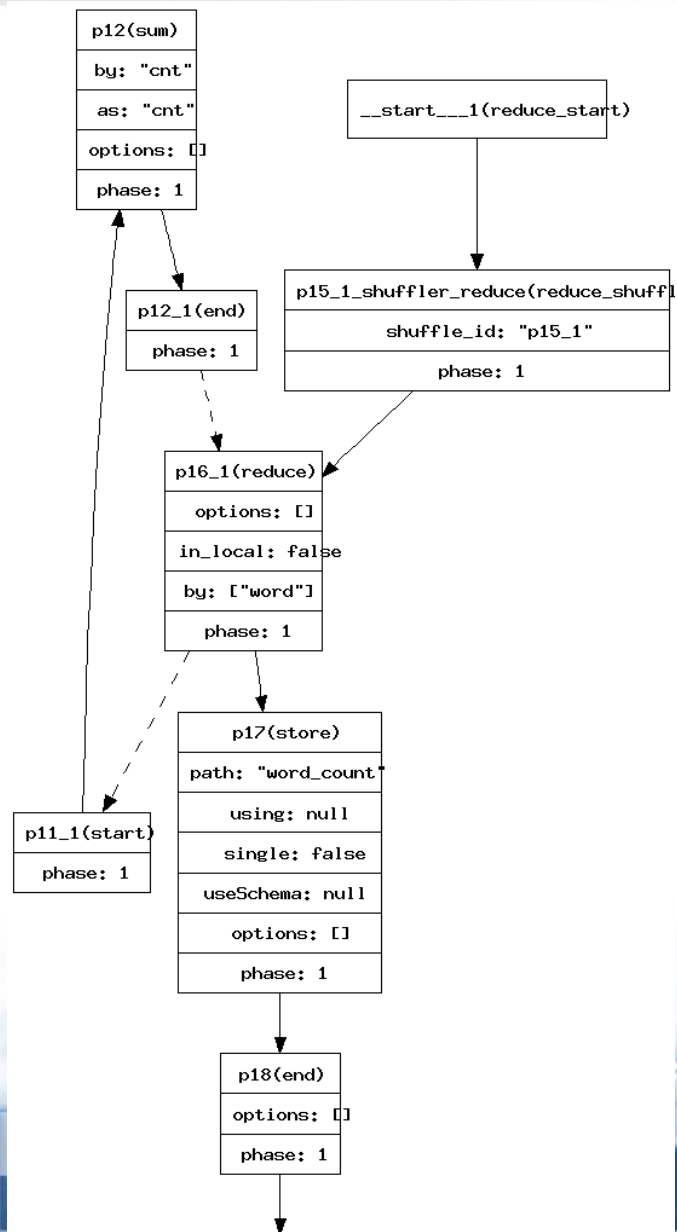
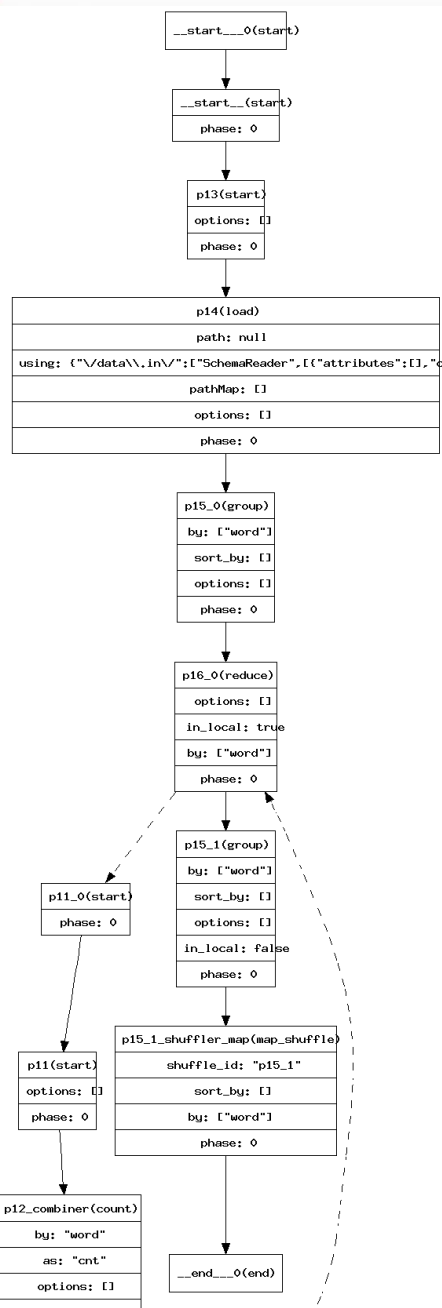


Field	ID	name	age	Score
Type	UInt64	string	int32	double
Index	2	5	9	10

- 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
  - 减少作业轮数、减少I/O、减少重复计算







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 ));
```

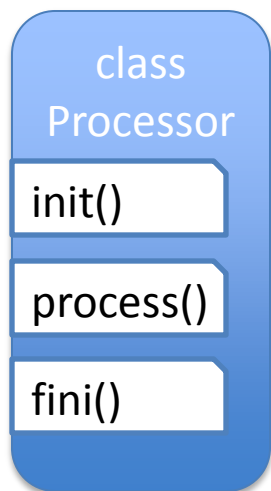
```
111 $p12_combiner->set_to('p12');
112
113 $p12_0 = new Ender;
114 $p12_0->set_id('p12_0');
115
116 $p15_1 = new Grouper(array
117     0 => 'word',
118 );
119 $p15_1->set_id('p15_1');
120
121 $p15_1_shuffler_map = new
122     0 => 'word',
123 ),NULL,array (
124 );
125 $p15_1_shuffler_map->set_id
126
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         p16_1->follow(*p15_1_shuffler_r
40         p17->follow(*p16_1);
41         p18->follow(*p17);
42         __end__->follow(*p18);
43         p15_1_shuffler_reduce->follow(*
44
45         __end__1->follow(*__end__);
46
47         /***** exec section *****/
48         bsl::var::IFactory& factory = C
49         Value main_record(factory);
50         __start__1->init(0);
51         __start__1->process(0, main_re
52         __start__1->fini(0);
53         /***** clean section *****/
54         delete p12;
55         delete p12_1;
56         delete p11_1;
57         delete p16_1;
58         delete p17;
59         delete p18;
60         delete __end__;
61         delete p15_1_shuffler_reduce;
62         delete __start__1;
63         delete __end__1;
64     } catch(bsl::Exception &e){
65         fprintf(stderr, "BSL EXCEPT
66     } catch(std::exception &e){
67         fprintf(stderr, "STD EXCEPT
68     } catch(...){
69         fprintf(stderr, "RUNTIME ER
```

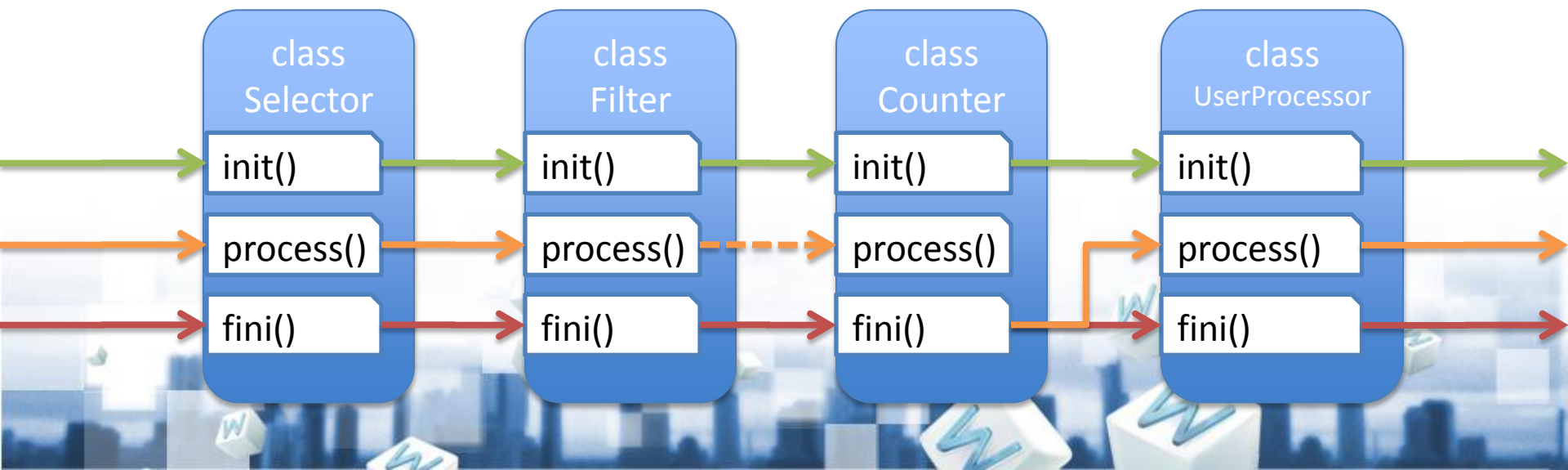
## Processor模型——Pipes & Filter模式



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



## ■ 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
简单编辑	24%
DQuery模式	43%
复杂编辑	33%

} 67%

## ■ LSP平台用户数量

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%

LOG中自有黄金屋

日志分析基本过程

百度日志分析成长历程

深入LSP平台

深入DISQL语言

总结与问答



- 日志的价值
  - 了解用户 • 了解自己
- 日志分析基本过程
  - 提取与传输 • 解析与过滤 • 计算 • 使用（报表、图表、回馈线上...）
- 百度日志分析成长历程
- 深入LSP平台
  - 平台架构 • 平台UI • 三种编辑模式
- 深入DISQL语言
  - 一个例子 • 几个特点 • 前端处理 • 中间语言翻译 • 运行时
- 采用情况
- 问题
  - 亦可通过微博([@陈晓鸣在百度](#))或邮箱([chenxiaoming@baidu.com](mailto:chenxiaoming@baidu.com))提问
  - 请关注Hadoop in China大会12月2日2时20分：《[DISQL2.0](#)》





我们还有很多**非常有挑战、非常有用、非常好玩**的问题急需解决.....

如果**你**恰好也喜欢解决这些问题.....  
请发邮件到[chenxiaoming@baidu.com](mailto:chenxiaoming@baidu.com)  
加入**百度**，让我们一起来把它们解决！

谢谢！

