

FEX Web 前端研发部

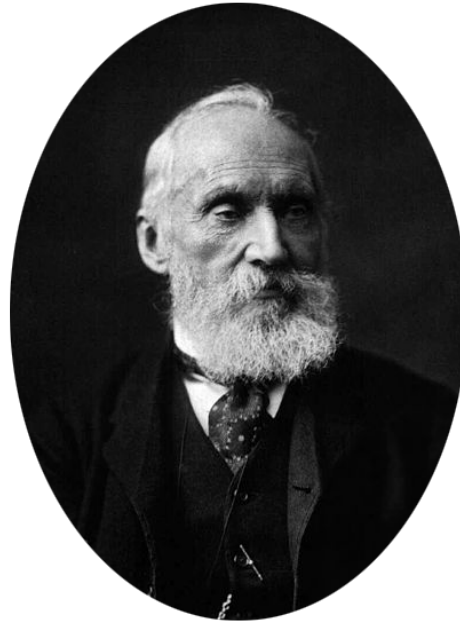
前端数据之美

张军-@大人物oo7

目录

- 一. 前端监控方式
- 二. 如何采集前端数据
- 三. FE也能玩大数据
- 四. 我们还能做什么
- 五. 总结

Part 1 前端监控模式



“If you cannot measure it, you
cannot improve it”

——William Thomson

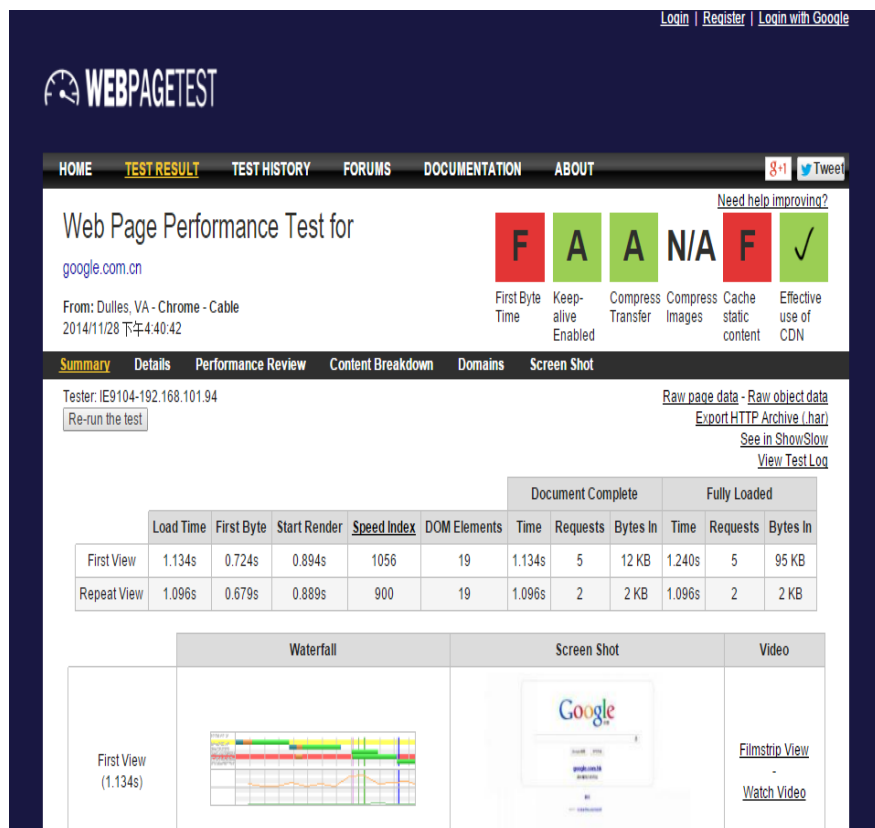


WE NEED YOU Data!!!

监控方式：非侵入式

WebPageTest

输入url获取各地测试数据



PhantomJS

命令行工具

The screenshot shows the PhantomJS website. At the top, there's a navigation bar with links like SOURCE CODE, DOCUMENTATION, API, EXAMPLES, and FAQ. Below this, the main heading is "Full web stack No browser required". A description states: "PhantomJS is a headless WebKit scriptable with a JavaScript API. It has fast and native support for various web standards: DOM handling, CSS selector, JSON, Canvas, and SVG." Below this, there are two buttons: "Download v1.9" and "Get started". On the right side, there's a code block showing a simple JavaScript example.

```
Simple Javascript example

console.log('Loading a web page');
var page = require('webpage').create();
var url = 'http://www.phantomjs.org/';
page.open(url, function (status) {
    //Page is loaded!
    phantom.exit();
});
```

At the bottom, there's a "Community" section with links to "Read the release notes", "Join the mailing list", and "Report bugs".

PhantomJS is an optimal solution for

HEADLESS WEBSITE TESTING

Run functional tests with frameworks such as Jasmine, QUnit, Mocha, Capybara, WebDriver, and many others.

SCREEN CAPTURE

Programmatically capture web contents, including SVG and Canvas. Create web site screenshots with thumbnail

PAGE AUTOMATION

Access and manipulate webpages with the standard DOM API, or with usual libraries like jQuery. [Learn more](#)

NETWORK MONITORING

Monitor page loading and export as standard HAR files. Automate performance analysis using YSlow and Jenkins. [Learn more](#)

监控方式：侵入式

JS监测

页面中嵌入监控脚本

Google Google Analytics (分析)

搜索此网站



首页 功能 学习资源 合作伙伴 帮助

登录 或 创建帐户

企业级的网站分析服务。

基于 Google 的世界级平台。 [了解详情](#)

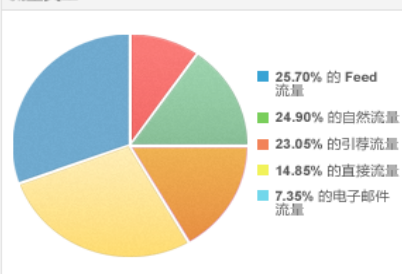
首页 标准报告 自定义报告

我的信息中心

每日访问次数



流量类型



不同国家/地区的网站停留时间

国家/地区	访问次数	平均网站停留时间
美国	67,445	00:01:54
英国	18,948	00:01:37
印度	8,882	00:00:58
加拿大	6,371	00:01:02
德国	5,845	00:00:32
法国	5,243	00:00:38

监控方式对比

类型	优点	缺点	示例
非侵入式	<ul style="list-style-type: none">指标齐全客户端主动监测竞品监控	<ul style="list-style-type: none">无法知道影响用户数采样少容易失真无法监控复杂应用与细分功能	Yslow , WebPageTest , PhantomJS , 百度UAQ
侵入式	<ul style="list-style-type: none">真实海量用户数据能监控复杂应用与业务功能用户交互、异步操作	<ul style="list-style-type: none">需插入脚本统计浏览器兼容性无法监控竞品	Google统计 , 各种用户行为监控

我们的做法

- 使用JS监测线上用户真实访问性能 为主
- 使用PhantomJS等工具线下分析 为辅

Part 2 如何采集前端数据

一、监控原理



二、监控框架

1. 统一发送数据
2. 配置抽样命中监控模块
3. 异步加载监控模块
4. 提供公共的on、un、fire等事件操作

Alog : <https://github.com/fex-team/alog>

三、监控模块

- ① 访问痕迹
- ② 性能数据
- ③ JS异常
- ④ 其他：点击分析、跨域资源统计等

① 访问痕迹

```
// 创建pv监控模块
alog('pv.create', {
  postUrl: 'http://localhost:8080/pv.gif',
  pageId: 'tieba-index',
  refer: document.referrer
});
// 发送数据
alog('pv.send', 'pageview');
```



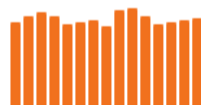
<http://localhost:8080/pv.gif?pageId=tieba-index&type=pageview&refer=baidu.com>

数据展示

PV 750,159,367

同比变化 +9.37% (2014-03-13)

PV波动



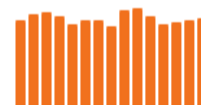
(2014-03-06~2014-03-20)

UV

643,459,874

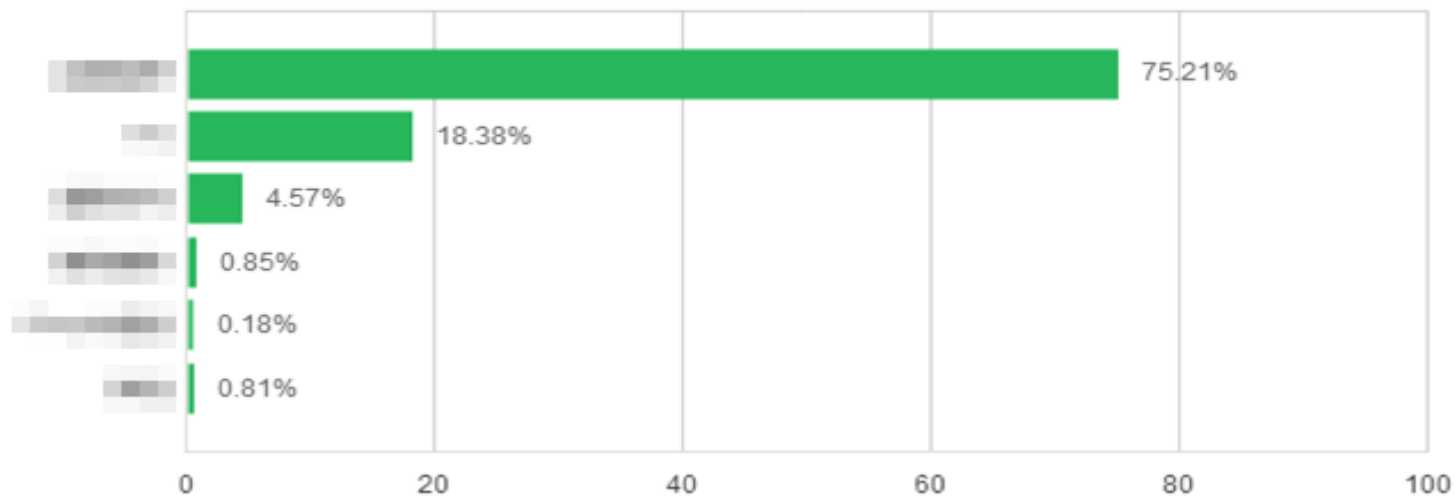
同比变化 +8.10% (2014-03-13)

UV波动

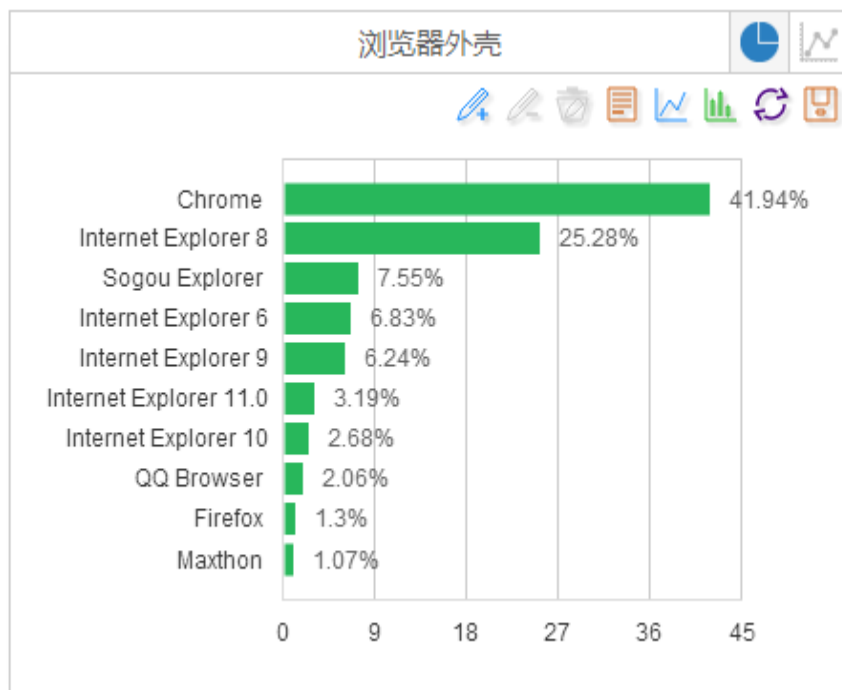


(2014-03-06~2014-03-20)

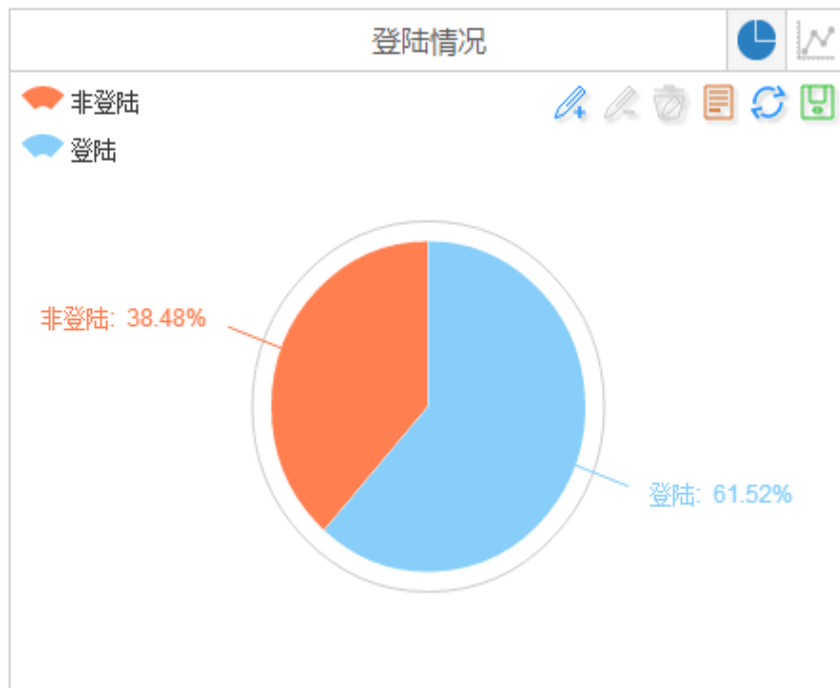
来源分布



数据展示

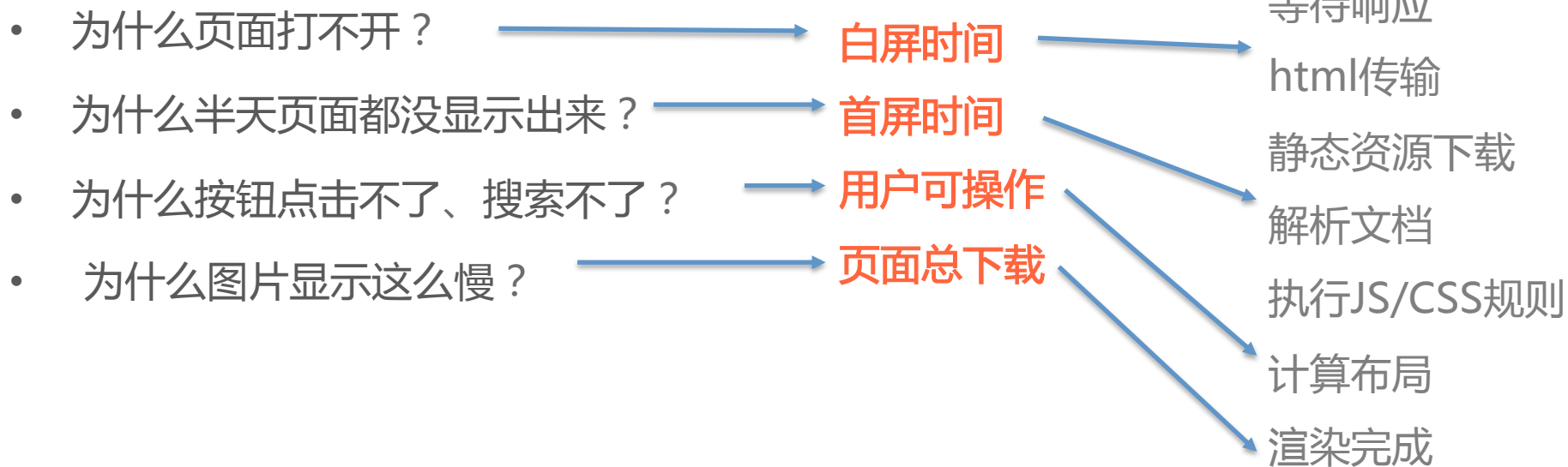


数据展示



② 性能数据

确定核心指标



时间起点

用户点击一个超链接或者输入url确认开始统计

方式一

使用cookie/hash记录用户点击超链接的时间戳

方式二

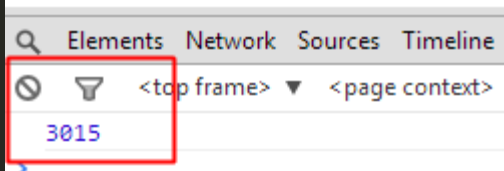
使用 [Navigation Timing API](#)

白屏时间

头部资源下载完成 \approx 白屏时间

```
<meta charset="UTF-8"/>
<script>
    var start_time = +new Date; //测试时间起点，实际统计起点为DNS查询
</script>
<!-- 3s后这个js才会返回 -->
<script src="script.php"></script>
<script>
    var end_time = +new Date; //时间终点
    var headtime = end_time - start_time; //头部资源加载时间
    console.log(headtime);
</script>
</head>
<body>
    <p>在头部资源加载完之前页面将是白屏</p>
    <p>script.php被模拟设置3s后返回，head底部内嵌JS等待前面js返回后才执行</p>
    <p>script.php替换成一个执行长时间循环的js效果也一样</p>
</body>
```

script.php替换成一个执行长时间循环



head底部内嵌JS来统计头部资源加载 -> 白屏时间

首屏时间



图片是制约首屏的主要因素

获取首屏内图片的加载耗时即可获取首屏时间

首屏时间—统计流程



3 页面onload之后找到最慢一张图片加载时间

1 首屏大概位置嵌入统计JS

首屏时间—陷阱

1. 图片加载完成、出错，gif图片重复触发加载事件的处理
2. iframe的处理：同图片
3. 异步渲染的处理：异步数据插入后再计算首屏
4. css背景图片的处理：首屏重要css背景通过JS发起图片请求判断是否已加载
5. 没有图片则以文字出现时间为准，可认为此统计首屏内JS执行时间

用户可操作时间

DomReady or 核心模块JS加载完毕(JS中打点标记)

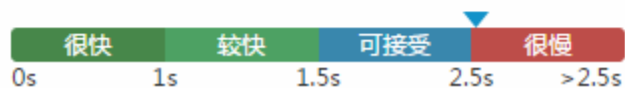
总下载时间

onload or 异步渲染完成

数据展示

首屏时间(昨日) ^③

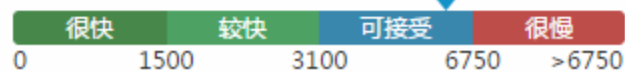
2.6s 很慢



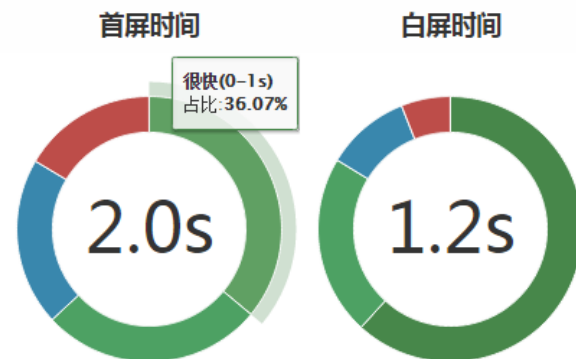
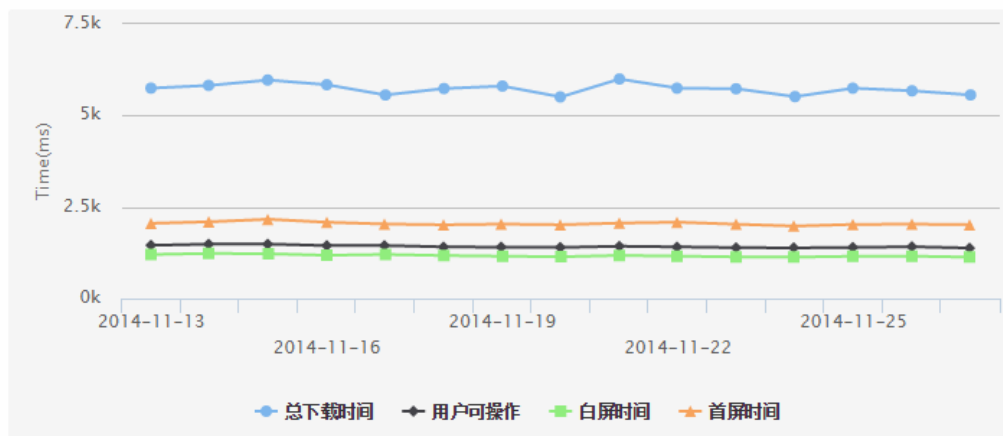
页面综合指数(昨日) ^③

6027 用户可接受

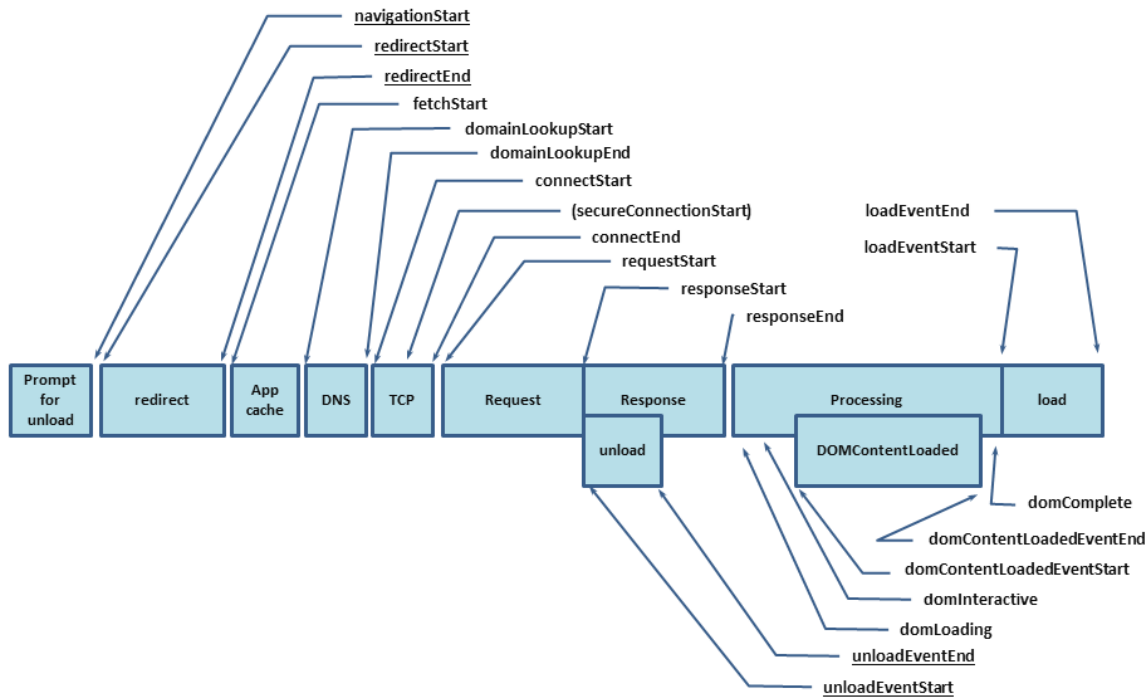
↓1.4% 正常波动



数据展示



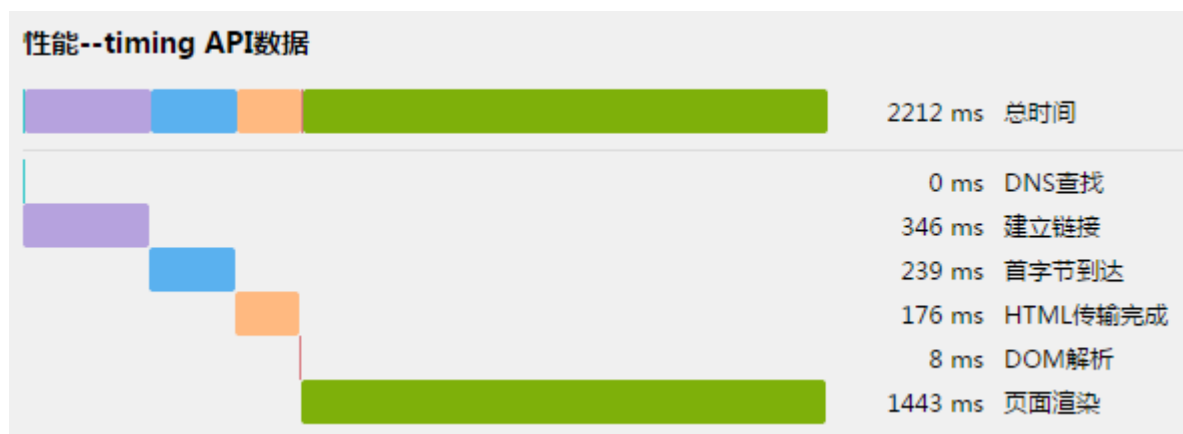
Navigation Timing API



```
> window.performance.timing
▼ PerformanceTiming {LoadEventEnd: 13884033503
1388403350128...} ⓘ
  connectEnd: 1388403349694
  connectStart: 1388403349667
  domComplete: 1388403350296
  domContentLoadedEventEnd: 1388403350139
  domContentLoadedEventStart: 1388403350128
  domInteractive: 1388403350128
  domLoading: 1388403349759
  domainLookupEnd: 1388403349667
  domainLookupStart: 1388403349663
  fetchStart: 1388403349658
  loadEventEnd: 1388403350311
  loadEventStart: 1388403350310
  navigationStart: 1388403349617
  redirectEnd: 0
  redirectStart: 0
  requestStart: 1388403349694
  responseEnd: 1388403349765
  responseStart: 1388403349744
  secureConnectionStart: 0
  unloadEventEnd: 0
  unloadEventStart: 0
  ▶ proto : PerformanceTiming
```

<http://www.w3.org/TR/navigation-timing/>

数据展示



Chrome插件:

<https://chrome.google.com/webstore/detail/web-performance/dmndepgifbfcjjmkjfgoecldegiodnjc>

③ JS异常

初期： `window.onerror`

优点

- 可以监控到几乎所有JS异常
- 产品线的JS代码无需任何修改

缺点

- 线上JS是混淆压缩的，无行号
- 跨域的JS异常，错误信息是 “Script error”
- 跨域的JS异常，获取不到JS文件名

产品线JS的CDN基本都跨域



Script error

PV异常率：0.46%(648,147) UV异常率：0.88%(544,099) 异常发生次数：5,650,925

异常发生的文件列表：

(1) 未知文件 PV异常率：0.43%(636,688) UV异常率：0.17%(533,301) 异常发生次数：437,837

③ JS异常

Then : try/catch

优点

- 不怕跨域
- 无惧JS压缩和混淆，捕获压缩前的行号
- 精确定位到模块、方法名

难道要手动给所有JS加try/catch？

```
a.js
1  var a = {
2      init: function(){
3          try{
4              hello();
5          }catch(e){
6              alog('exception.send',{
7                  msg: e.message,
8                  line: 2
9                  module: 'a.js',
10                 method: 'init'
11             });
12         }
13     }
14 }
```

③ JS异常

我们的做法

- 前端使用FIS开发 (fis.baidu.com)
- 使用我们提供的FIS插件 (<https://github.com/zhangjunah/fis-preprocessor-js-exception-catch>)
- JS文件打包时，自动在function内部加上try/catch

```
//编译前
var slide = {
  init: function(name){
    this.name = name;
  },
  name: ''
}

//编译后
var slide = {
  init: function(name){try{
    this.name = name;
  }catch(e){.....}},
  name: ''
}
```

```
//编译前
$('#search').on('click', function(){
  search();
});

//编译后
$('#search').on('click', function(){try{
  search();
}catch(e){.....}});
```

数据展示

PV异常率 0.12%(557,876)

同比变化 +0.19% (2014-02-24)

PV异常率波动



(2014-02-17~2014-03-03)

异常PV浏览器



Safari Chrome IE8 QQ Browser IE6 Sogou Explorer IE7 IE9
Other IE11 Firefox IE10 Opera

error test

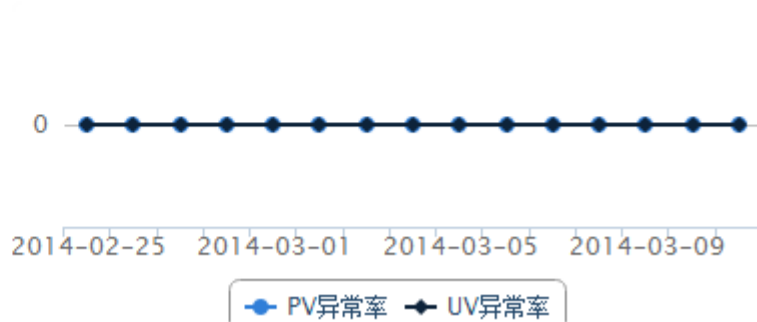
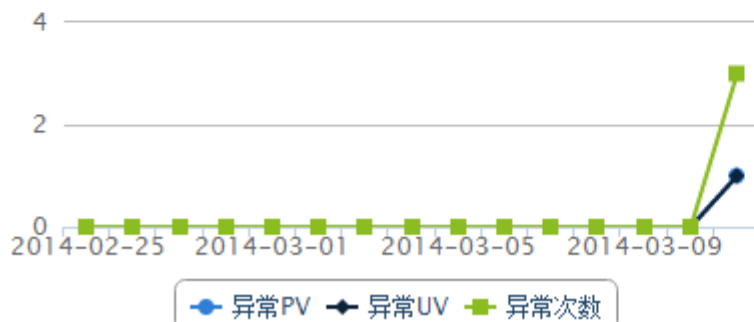
PV异常率: 0.00%(1) UV异常率: 0.00%(1) 异常发生次数: 3

异常发生的位置列表:

(1) common/export.js 异常method: call 所在行: 7 PV异常率: 0.00%(1) UV异常率: 0.00%(1) 异常发生次数: 3

异常发生的浏览器: Chrome(100%)

最近半月趋势:



④ 其他监控模块

点击数据

人均点击 **10.93**

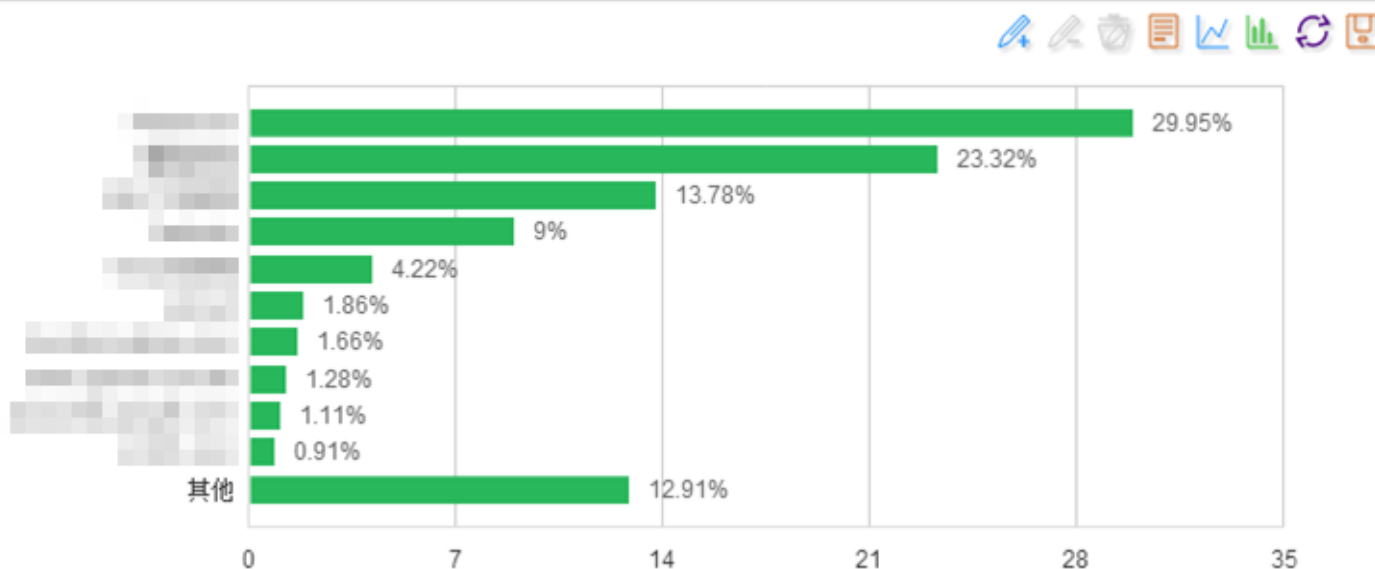
同比变化 **+7.82%** (2014-03-13)

人均点击波动



(2014-03-06~2014-03-20)

流出url分布



④ 其他监控模块

跨站资源

http://re.taotaosou.com/js/_tts_browser_center.js

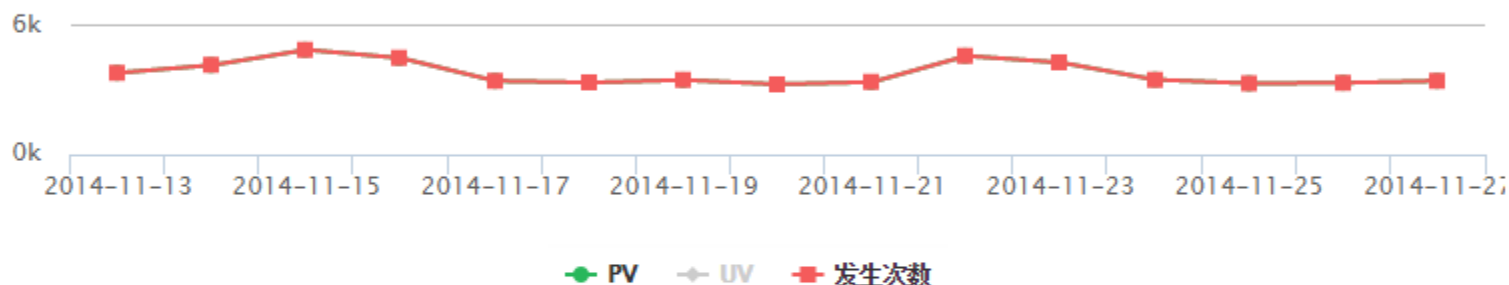
淘淘搜比价

修改备注

PV : 3,444 UV : 3,258 发生次数 : 3,446

操作: [查看详细日志](#)

最近半月走势



<http://202.102.100.100/35ff706fd57d11c141cdefcd58d6562b.js>

添加备注

PV : 3,036 UV : 2,892 发生次数 : 3,036

<http://re.taotaosou.com/browser-static/tmt/sea.js>

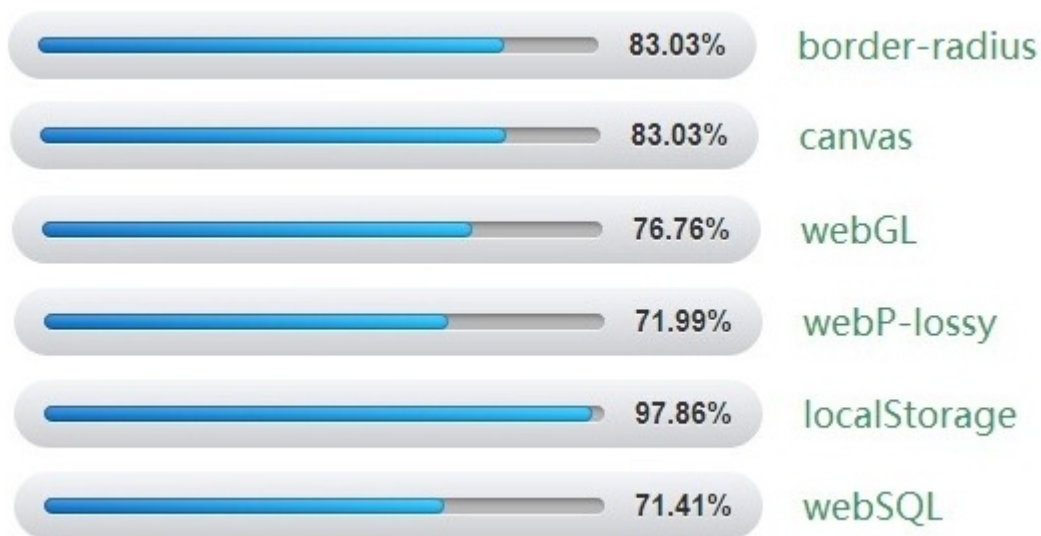
淘淘搜比价

修改备注

PV : 2,581 UV : 2,443 发生次数 : 2,581

④ 其他监控模块

HTML5、CSS3支持程度



<http://modernizr.com/>

本节回顾

1. 统计框架 Alog
2. 监控访问痕迹
3. 性能数据：白屏、首屏、用户可操作、总下载、timing API
4. JS异常捕获：onerror、try/catch、FIS方案
5. 其他监控模块：点击、跨站资源、H5/CSS3支持度

Part 3 FE也能玩大数据



- map/reduce
- <http://hadoop.apache.org>



- 类SQL语句来执行hadoop任务
- <http://hive.apache.org>

nodeJS + hadoop

Hadoop支持标准输入输出，nodeJS也同样支持

```
1  hadoop-szwg-stoff/bin/hadoop jar hadoop-szwg-stoff/contrib/streaming/*streaming*.jar \  
2  -D mapred.map.tasks=10 \  
3  -D mapred.reduce.tasks=1 \  
4  -D mapred.job.name=test001 \  
5  -cacheArchive 'hdfs://szwg-stoff-hdfs.dmop.baidu.com:54310/node-v0.8.4-linux-x64.tar.gz#nodejs' \  
6  -file map.js \  
7  -file reduce.js \  
8  -mapper map.js \  
9  -reducer reduce.js \  
10 -input test001.log \  
11 -output result
```

nodeJS + hadoop

在map.js、reduce.js中

```
// 表明运行环境为nodejs
#!/nodejs/bin/node

// 标准输入输出
var stdin = process.stdin;
var stdout = process.stdout;

// 处理每行数据
function eachLine(line) {
    // 这里以pv数据为例
    if (line.type == 'pageview') {
        // tieba-index \t pageview
        stdout.write(line.pageId + '\t' + 'pageview' + '\n');
    }
}

// 读数据流
stdin.on('data', function(chunk) {
    // 将chunk拆为每行
    for ( ..... ) {
        eachLine(line);
    }
});

// 读数据流结束
stdin.on('end', function() {
    console.log('hi hi, read end!!!');
});
```


Hive

创建数据表、导入hadoop格式化后的数据

```
CREATE TABLE data_all (  
    pageId STRING,  
    type STRING  
)  
ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' STORED AS TEXTFILE;  
  
LOAD DATA INPATH '/result' OVERWRITE INTO TABLE data_all;
```

查询所有页面的pv

```
SELECT pageId, count(*) AS pv  
FROM data_all  
GROUP BY pageId
```

Part 4 我们还能做什么

非侵入式监控

竞品对比分析、全互联网网页性能排名、命令行工具 等

Site.Archive

性能趋势

对比分析

详情分析

性能排名

使用帮助



请输入url...

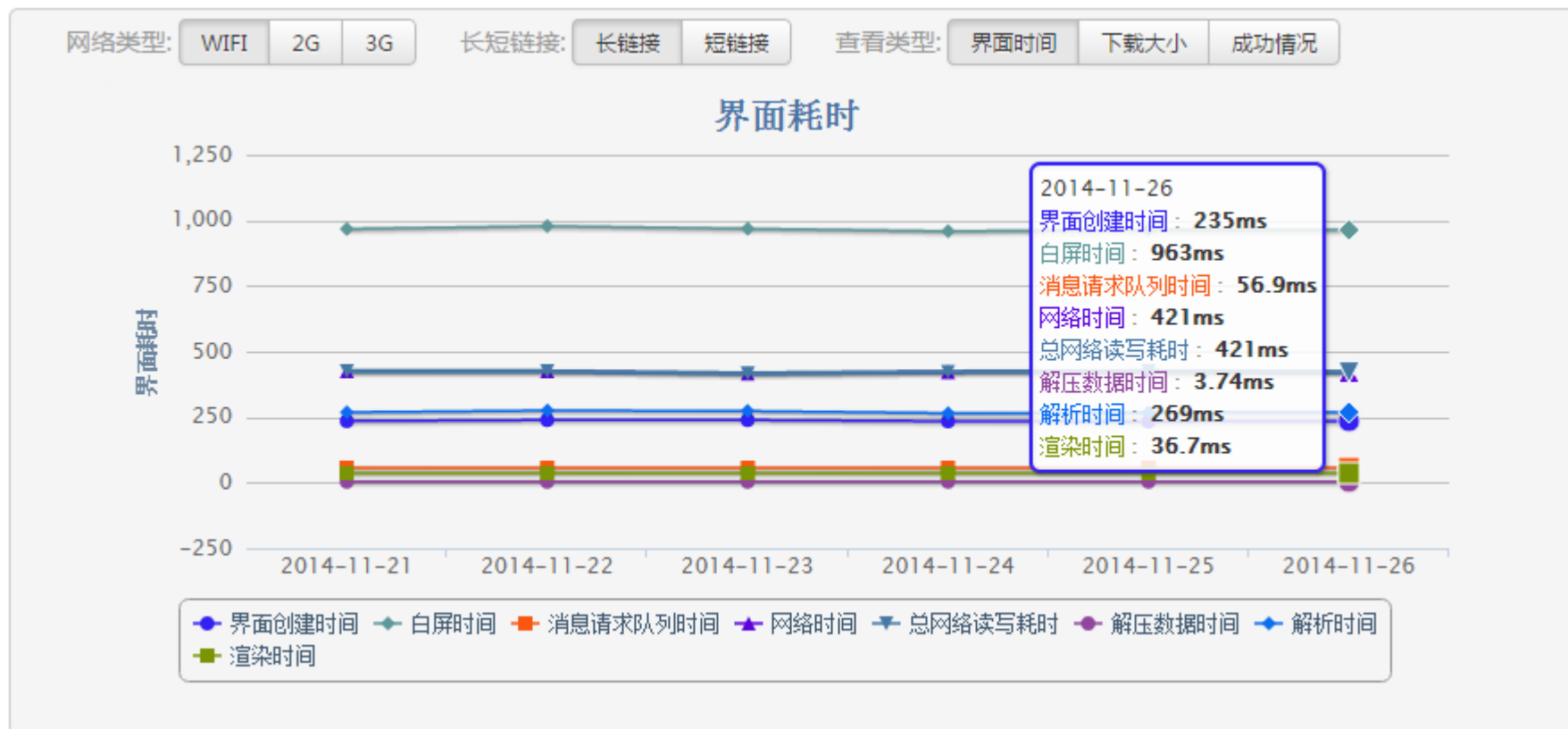
提交url

从全国上千网站采集每日数据，
汇聚加工成海量的数据库，为分析
和展现提供了坚实的数据支持，让
你看到最真实的互联网大数据下的
资源分配。

海量的数据支持

<http://siteArchive.baidu.com>

native监控



前后端监控结合

1. 前、后端性能统一监控，形成完整的时间线
2. 前、后端异常统一监控，快速定位bug

Part 5 总结

1. 监控方式：非侵入式、侵入式
2. 数据采集：Alog框架、访问、性能、JS异常等
3. 玩转大数据：nodeJS + hadoop 、hive
4. 下一步方向：非侵入式工具、native APP、前后端监控结合

Q&A

THANKS