# Semantic Matching in App Search

Juchao Zhuo, Zeqian Huang, Yunfeng Liu,
Zhanhui Kang, Xun Cao, Mingzhi Li, Long Jin

WSDM Feb 4, 2015

**TENCENT 腾讯**

# Outline

- ✓ Overview
- ✓ Challenge in App Search
- ✓ Semantic Matching
  - ● Matching with Topic
  - ● Matching with Tag
  - ● Learning to rank
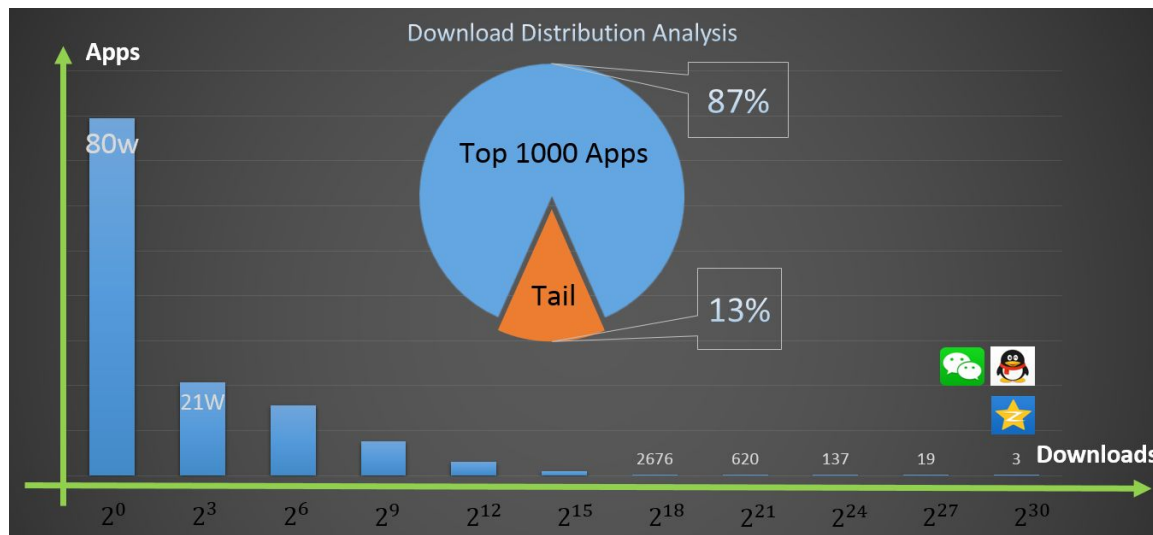- ✓ Applications and Evaluations
- ✓ Conclusion

# Overview

- ✓ With the rapid growth of smartphones, app market has become a significant mobile internet portal. As an important function in app market, app search gains lots of attentions.

- ✓ Miss-match is the critical challenge in app search. Semantic matching is a key technology to reduce miss-match.

- ✓ In this talk, we will describe a semantic matching platform , which mines topics and tags in big data to enrich query and app representations, and implements learning to rank.

- ✓ The semantic matching platform is used by "Myapp" app market, one of the top three android app markets in China.

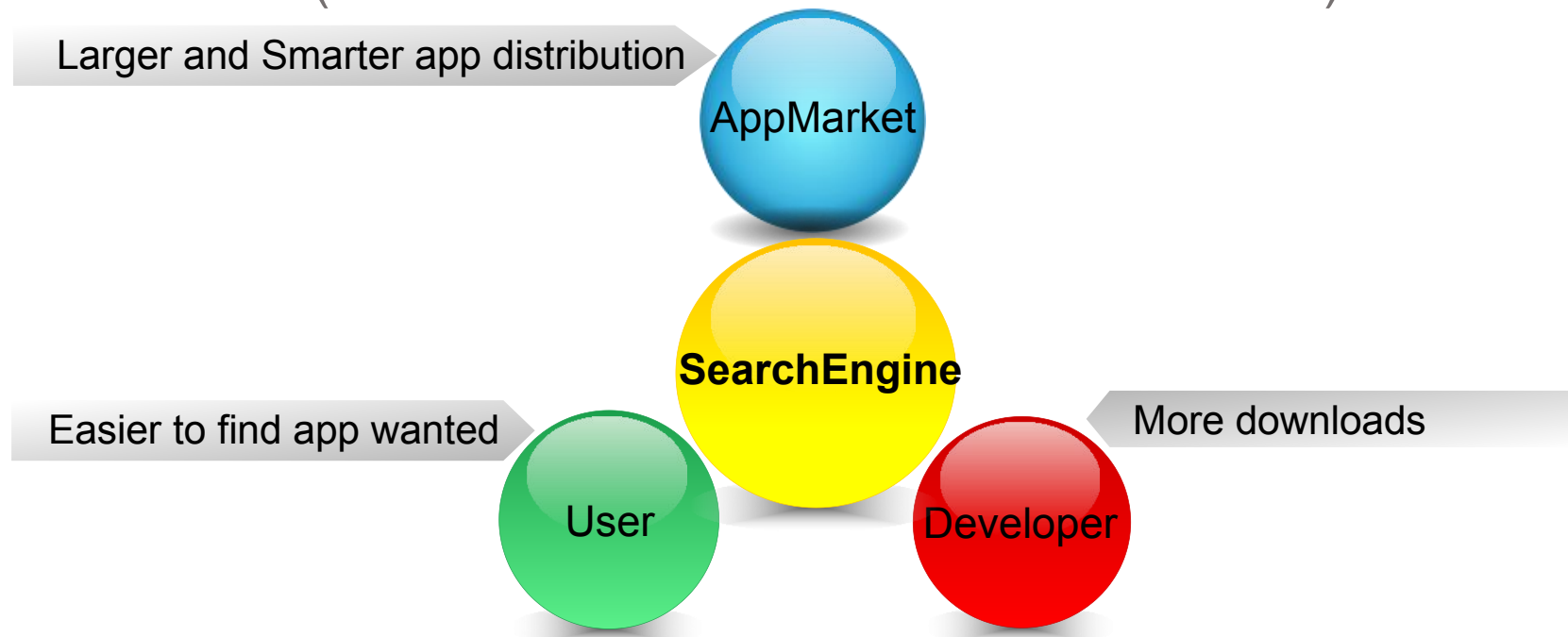# App Search Application in Tencent

✓ **MyApp**

- An android app market with a peak distribution of over 100 millions in one day of 2014
- App search engine contributes more than 40% to mobile app new-downloads
- Rapid growth: available apps from 0.3 million to 1.2 million within one year
- Long Tail: apps which were downloaded over a million times accounted for less than 0.1%

# Objectives of App Search

✓ App search objectives
- Facilitate the app market, users and developers

✓ App search metrics
- Downloads / QV(Query views) / UV(User views)
- CTR(Click-Through-Rate) / ROP(Rate of Penetration)
- NDCG (Normalized Discounted Cumulative Gain)

Larger and Smarter app distribution

AppMarket

SearchEngine

Easier to find app wanted

More downloads

User

Developer

5

# User habits in App Search

✓ Two kinds of queries in app search

| | Ratio of Query Number | Ratio of QV Number |
|---|---|---|
| Precise Search | 88% | 75% |
| Fuzzy Search | 12% | 25% |

✓ Precise Search
- Search by app name, mostly prompted by the search box

✓ Fuzzy Search
- Non-Name, always colloquial expression
- Content/Category/Function related
- User-habit of web search is brought to app search on mobile
- e.g.
  "微信里的游戏"(game in wechat)
  "音乐软件"(music application)
  "报时间的软件"(application that reminds time)

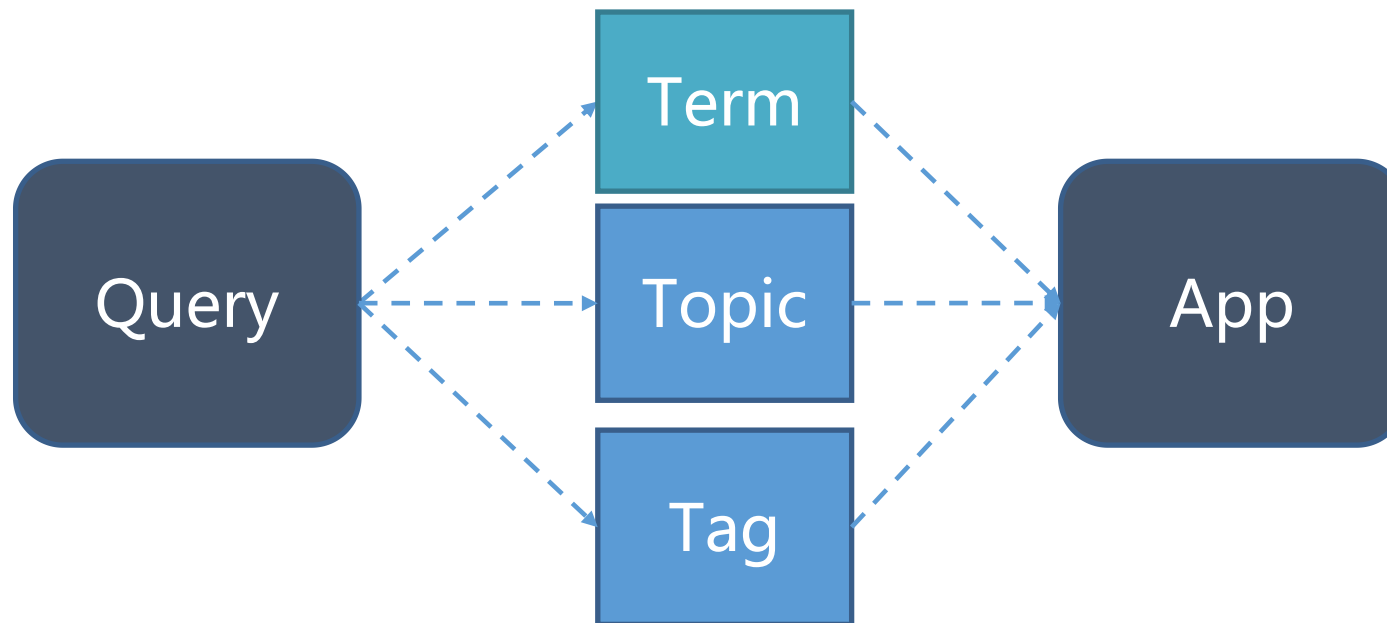# Challenge in App Search

✓ Miss-match

- ●Long tail challenge is more obvious in app search than that in web search.
- ●**Data shortage:** app data is much less than web data
- ●**Text shortage:** app name and desc. is the only text annotation for app
- ●Miss-match occurs when user and developer use different terms to describe the same semantic, traditional term matching can not fix it.
- ●e.g.



7

# Semantic Matching Methodology

✓ How to describe "Semantic "?
  ● Term + Topic + Tag
✓ Enrich query and app representations by topic and tag
✓ Perform query-app matching with the representations
✓ Hybrid Ranking Model: LTR

# Matching with Topic

✓ **Topic Model**

- Using Layered LDA(Latent Dirichlet Allocation) model
- MPI based parallel computing framework
- Topic probability distribution over term space: P(word|topic)



- Assign million apps to **1000+ topics**
- Doc probability distribution over topic space: P(topic|doc)



| TopicId | Category | Probability | TopicWords |
|---|---|---|---|
| 31 | 跑酷[游戏] | 0.625815 | 游戏,跳跃,障碍,收集,跑酷,躲避,动作,不断,控制,路上,奔跑,道具,避开,挑战,逃亡, |
| 26 | 关卡[游戏] | 0.0649123 | 游戏,玩家,关卡,挑战,难度,益智,休闲,玩法,画面,道具,过关,闯关,乐趣,等级,一共, |
| 70 | rpg[游戏] | 0.0593985 | 游戏,玩家,系统,技能,战斗,体验,丰富,特色,装备,角色,画面,玩法,华丽,全新,网游, |

# Matching with Topic

✓ Query inference with topic
- Each query is regarded as a document
- **Challenge:**short text has not enough information to inference
- **Solution:** Expanding query with collection of click apps

✓ Topic matching
- Map query and documents in the topic space
- Query $Vq$ and documents $V_d$ are both represented with probability distributions over topics
- Calculate topic match score between $Vq$ and $V_d$

query          document(app)

Topic space

$Vq$                $V_d$

# Matching with Tag

✓ **Limitation of Topic Matching**

- Text corpora of app documents is not large enough to support large topic number.

- Significant difference may still exist between apps in the same topic.

- Long tail queries lack statistical click data to expand, even after topic matching, many tail queries are still unknown.
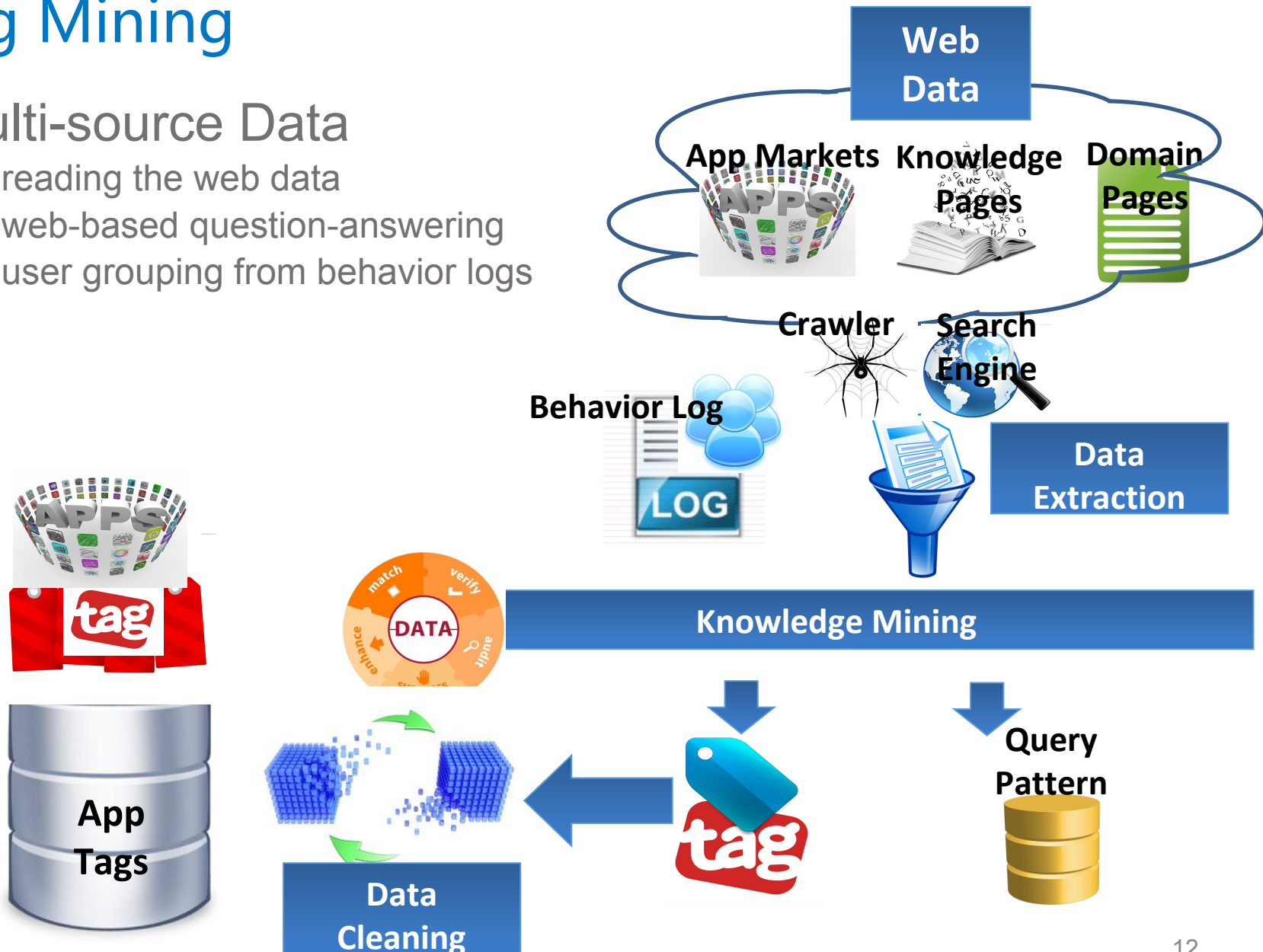
✓ **Matching with Tag**

- More fine-grained semantics can be described by tag.

- Most app stores assemble tags on human editorial curation.

- Our system can monitor the app ecosystem in real time, and automatically extract tags and assign them to apps from multi-source data.

# Tag Mining

✓ Multi-source Data
  ● reading the web data
  ● web-based question-answering
  ● user grouping from behavior logs

**Web Data**

**App Markets**   **Knowledge Pages**   **Domain Pages**

APPS

**Crawler**   **Search Engine**

**Behavior Log**

LOG

**Data Extraction**

**Knowledge Mining**

match  verify
**DATA**
enhance  audit

**App Tags**

**Data Cleaning**

tag

**Query Pattern**

12

# Tag extraction from web data

✓Data from web

- Structured data

| 中文名 | 全民枪战 APP | 发行日期 | 2014年 |
| 其他名称 | CA | 音 乐 | 小旭音乐 |
| 游戏类型 | 射击 | 内容主题 | 战争 |
| 游戏平台 | ios，Android | 玩家人数 | 中型射击手游 |
| 发行商 | cmge中国手游 | 游戏下载 | 4399游戏盒、摸摸、小皮游戏 |

TAG

- Unsructured text
  - using template to extraction

HED

ATT

ADV    RAD              ATT        ATT

Root    最    疯狂    的    枪战    手游    全民枪战
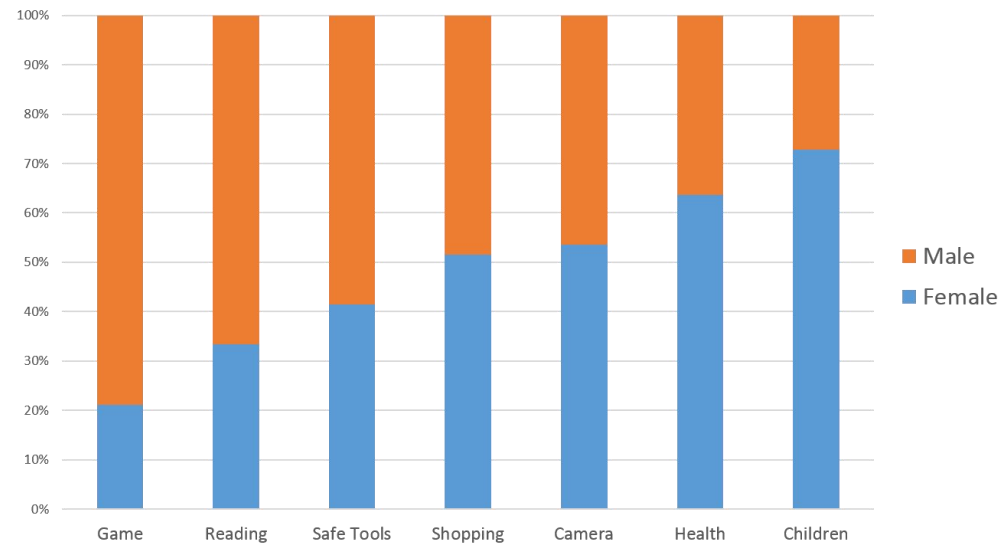                              TAG     TAG    APP ENTITY

# Tag from question-answering

● Using web-based question-answering to perform completion of missing tag-app pairs.

# Tag from user behavior logs

- Users Profile
  - gender, age, location, ...
- Users Behavior
  - search, download, install, ...

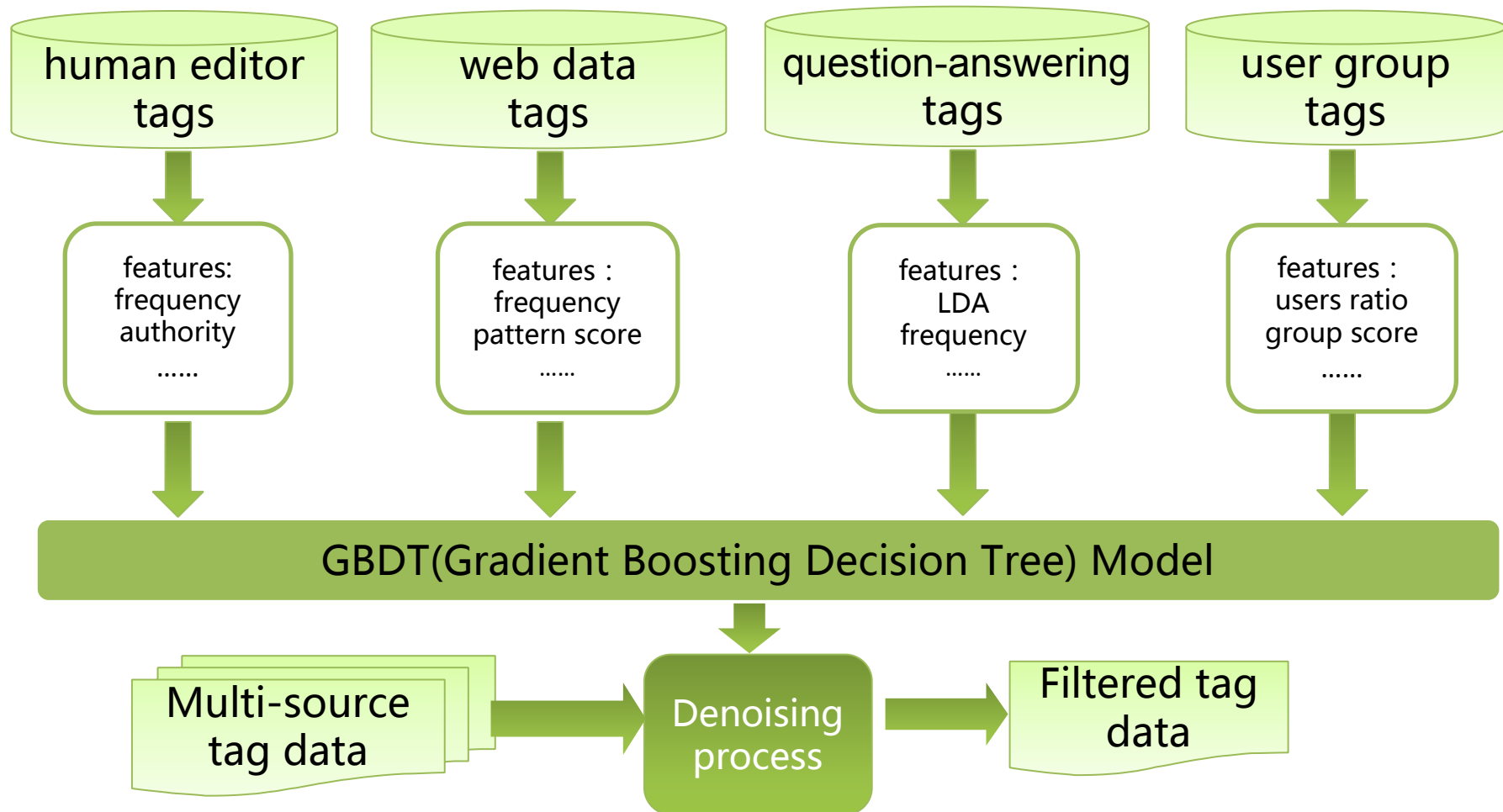Download-ratio by gender and category



- Based on user-behavior data and natural representation by tag

# Tag Denoising

- Using machine learning technology to calculate confidence

| human editor tags | web data tags | question-answering tags | user group tags |
|---|---|---|---|

| features: frequency authority …… | features : frequency pattern score …… | features : LDA frequency …… | features : users ratio group score …… |
|---|---|---|---|

**GBDT(Gradient Boosting Decision Tree) Model**

Multi-source tag data → Denoising process → Filtered tag data

# Tag statistic

✓Denoising data
  ●Over 97% pairs filtered: Treasures are hidden among the sands

| Denoising | |
|---|---|
| Original tag-app pairs | 38515130 |
| Filtered tag-app pairs | 37677471 (97.83%) |
| Valid tag-app pairs | 837659 (**2.17%**) |

✓Statistic data
  ●Over 90K tags mined
  ●Covers 83% of apps
  ●Top 100K apps have 8.53 tags in average

| Sources | Tags |
|---|---|
| Web Data | 81626 |
| Question-answering Data | 8620 |
| User Group | 1218 |
| Total | 91464 |

# Matching query with tag

✓Using template to map query to tags

query

Template match

Template Set

Template : Tag.+Category.+suffix

单机.+游戏.+哪种好
格斗.+游戏.+求推荐
赛车.+游戏.+下载

Q:单机游戏哪种好

Extract tag

T:单机+游戏

✓Using click data to calculate confidence of template
- P(Template) = 2/3

Q:单机游戏哪种好

Click Data

| APP:斗地主 | T:单机+游戏 | TAG retrieved |
| APP:找你妹 | T:单机+游戏 | TAG retrieved |
| APP:QQ游戏 | "游戏" | Term retrieved |

# Learning to Rank

✓ **Challenge of Ranking**
  - Relevance calculated by different matching model are incomparable
  - the example data is imbalanced (e.g. colloquial query less than normal)
  - Most of the features are nonlinear

✓ **Using LambdaMart to rank**
  - LambdaMART combines MART and LambdaRank to solve the supervised learning problem
  - Mart(Multiple Additive Regression Trees) is a gradient boosting tree model
  - Label training data partitioned by query
  - Maximizing NDCG by learning relevance score through MART

# Learning to Rank Application

✓ LambdaMART of Combine Ranking
  - 50+ different features
  - 300,000+ pairwise traning data
  - 3000+ test samples

✓ Offline Experiment measurement
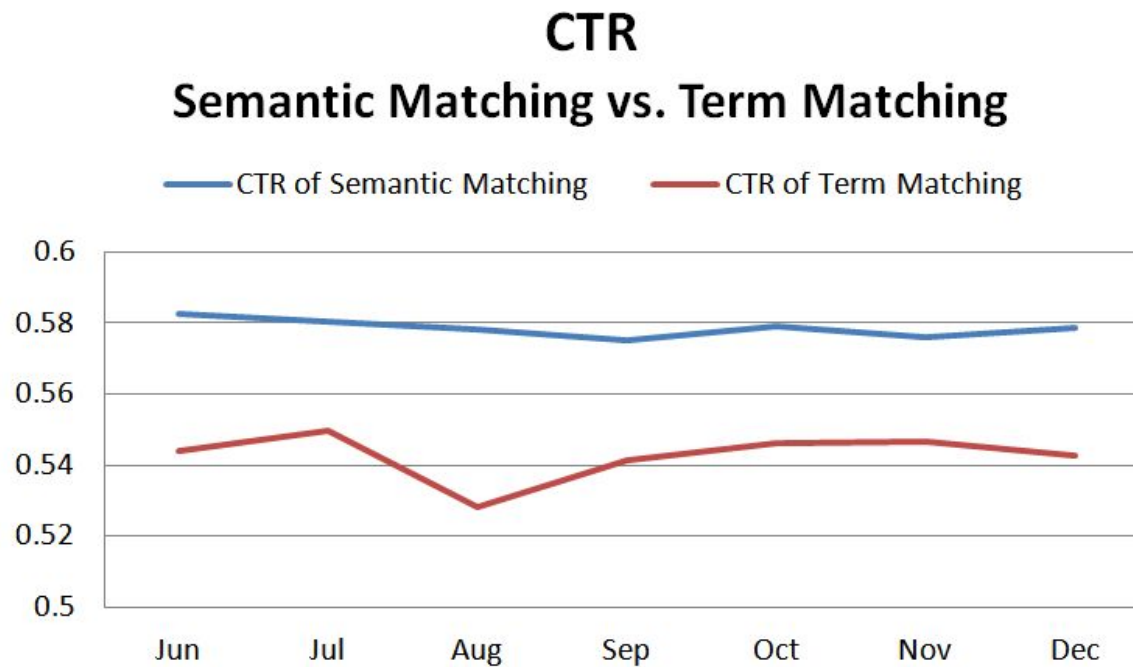
| NDCG of baseline | NDCG of LambdaMART | Improvement vs baseline |
|---|---|---|
| 0.8733 | 0.9553 | 9.4% |

✓ Online A/B Test measurement
  - CTR promoted by 6%↑

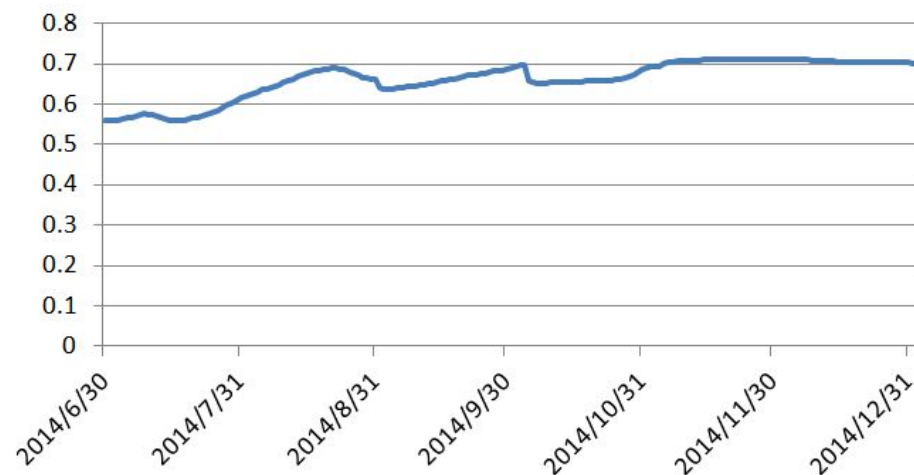# Semantic Matching Metrics

✓ Online A/B Test measurement
- CTR diff. on query samples
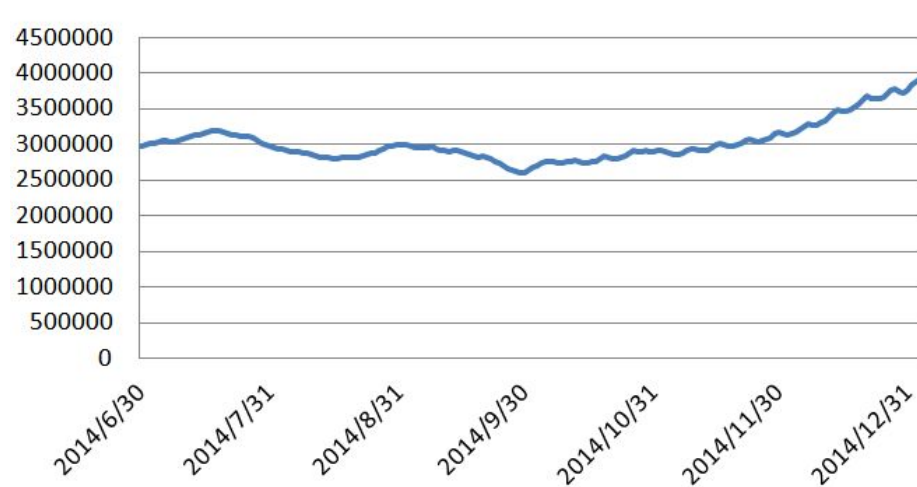- 9.7% querys & 26.9% query views
- CTR promoted by 6%~8%↑

## CTR
### Semantic Matching vs. Term Matching

— CTR of Semantic Matching    — CTR of Term Matching

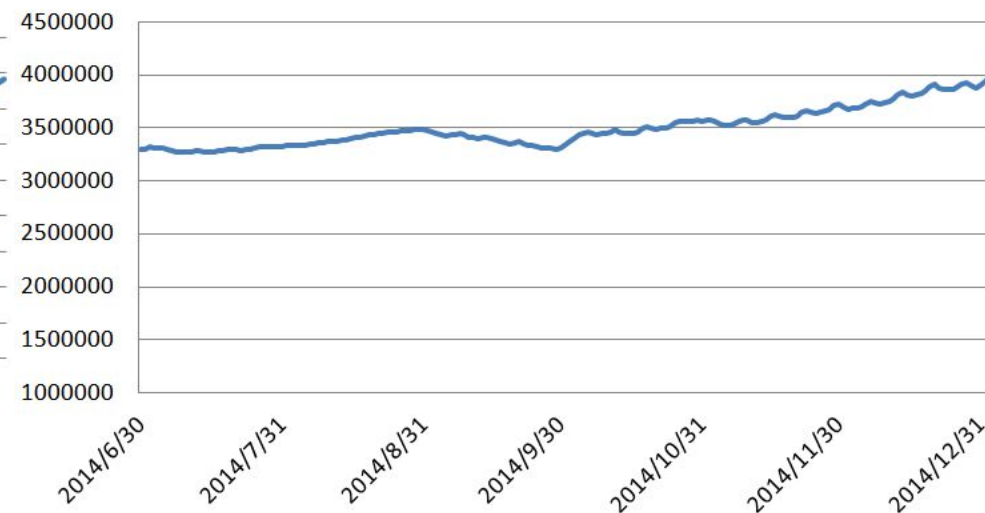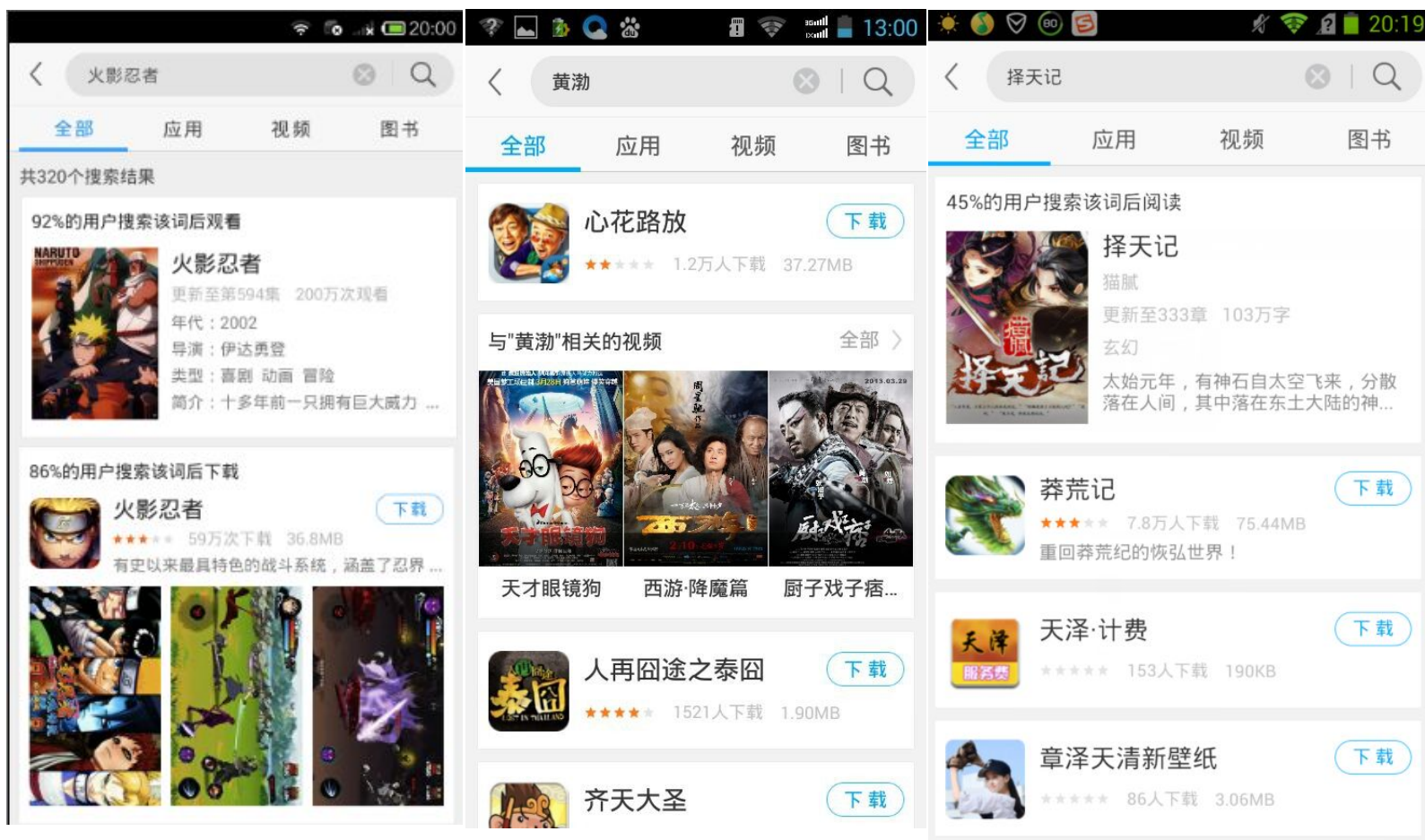# App Search with Semantic Matching

**CTR**



**Search UV**



**Download Apps**

# App Indexing

✓ Besides apps, all digital content inside apps can be offered

✓ Using LambdaMart to rank with different digital content

✓ ROP(Rate of Penetration) of App indexing version speed-up 32.3%↑

# Semantic Matching Application Example

✓ Deeper understanding colloquial form of query

# Conclusion

- ✓ From **Term matching** to **Semantic matching**
  - Richer representation of semantic

- ✓ Methodology of Enriching infomation
  - Use the web search technology to detect the relationship of app data

- ✓ What is the next direction of mobile search
  - More input mode: voice, photograph, two-dimension code
  - Search engine should become more intelligent

- ✓ Stay tuned for 2015!

# Acknowledgement

- ✓ Joint work with many brilliant colleagues from several departments of Tencent

- ✓ Many constructive inputs are from Yue Wu and Xing Yao

内部搜索平台部
Internal Search Platform Dept

腾讯应用宝
就 要 玩 在 一 起