

# AiCon

全球人工智能与机器学习技术大会

# iQIYI自然语言处理和视频大数据 分析应用

吴友政

Powered by 爱奇艺技术产品中心

主办方 **Geekbang** 极客邦科技 **InfoQ**

# 人工智能基础课

“通俗易懂的人工智能入门课”

王天一  
博士 副教授



扫一扫，免费试读

# AI技术内参

你的360度人工智能信息助理

洪亮劼  
Etsy 数据科学主管



扫一扫，免费试读



## 关注落地技术，探寻AI应用场景

- 14万AI领域垂直用户
- 8000+社群技术交流人员，不乏行业内顶级技术专家
- 每周一节干货技术分享课
- AI一线领军人物的访谈
- AI大会的专家干货演讲整理
- 《AI前线》月刊
- AI技能图谱
- 线下沙龙



扫码关注带你涨姿势

# QCon

## 全球软件开发大会

# 成为软件技术专家的 必经之路

### [北京站] 2018

会议：2018年4月20-22日 / 培训：2018年4月18-19日

北京·国际会议中心

**8折** 购票中, 每张立减1360元  
团购享受更多优惠



识别二维码了解更多

# ArchSummit

## 全球架构师峰会

2018 · 深圳站

从2012年开始算起，InfoQ已经举办了9场ArchSummit全球架构师峰会，有来自Microsoft、Google、Facebook、Twitter、LinkedIn、阿里巴巴、腾讯、百度等技术专家分享过他们的实践经验，至今累计已经为中国技术人奉上了近千场精彩演讲。

限时**7折**报名中，名额有限，速速报名吧！

● 2012.08.10-12 深圳站



● 2018.07.06-09 深圳站

会议：07.06-07.07

培训：07.08-07.09



# 个人介绍

- 教育背景
  - 中科院自动化所 ( NLPR , CASIA ) 毕业
- 工作经历
  - iQIYI
  - SONY China Research Lab.
  - University of Edinburgh , UK
  - NICT , Japan
- 参与项目/研究兴趣
  - 自然语言处理、语音助手、问答系统、机器翻译
  - 语音识别 ( Kaldi )
  - 商业智能



# TABLE OF CONTENTES

- 理解视频内容
  - 中文词法分析
  - 预测（票房和流量）
- 理解视频用户
  - 舆情监测
  - 查询理解和意图搜索
- 总结

# NLP

## 支持业务

泡泡

头条

VR

BI

广告

搜索

推荐

用户画像

客服中心

审核平台

## 应用研发

查询理解

打标签

语音助手

舆情监测

热点事件

## 中文词法分析

分词

词性标注

词权重

实体识别

实体链接

## 数据挖掘

知识图谱

领域词典

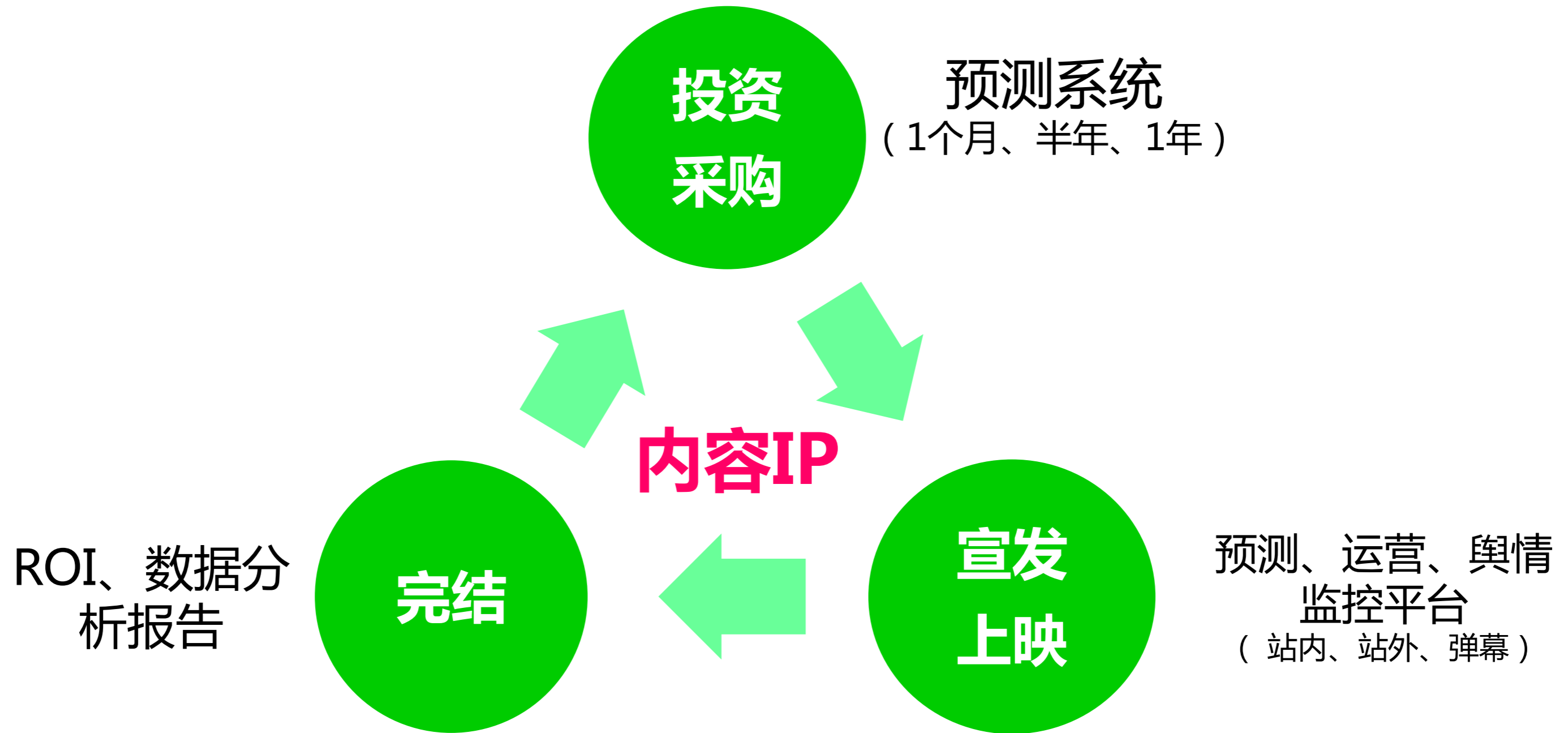
同义词

情感词

embedding



# 视频大数据分析



# TABLE OF CONTENTES

- 理解视频内容
  - 中文词法分析
  - 预测（票房和流量）
- 理解视频用户
  - 舆情监测
  - 查询理解和意图搜索
- 总结

# 中文词法分析

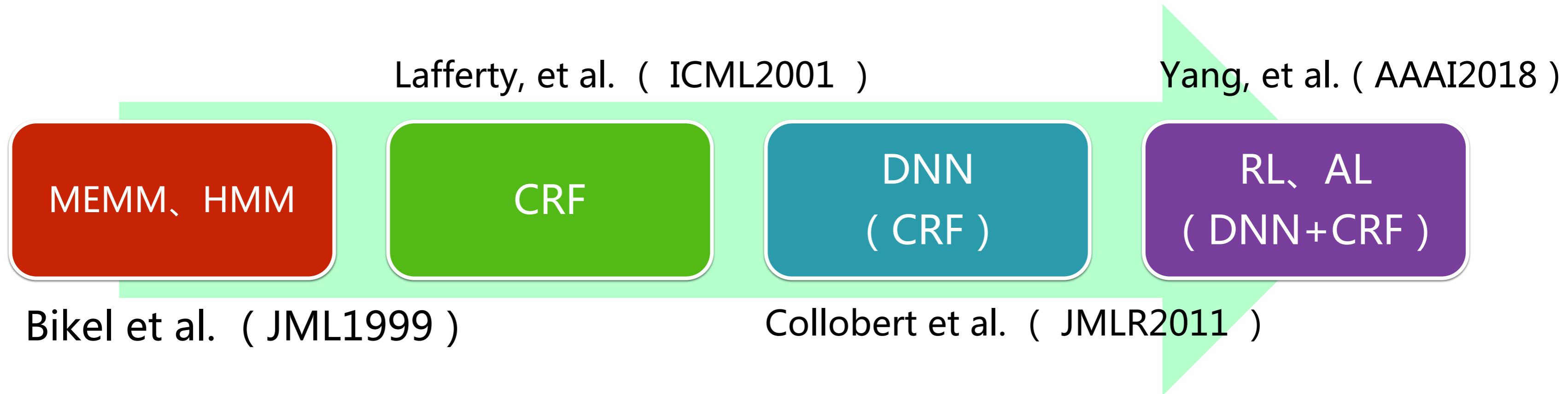


鹿晗	演绎	帅气	陈长生	逆天改命	择天记	终	迎	结局
person	v	adj	character	idiom	album	d	v	n
0.758	0.284	0.402	0.577	0.476	0.756	0.0	0.0	0.370

# 实体识别简介

- 识别文本中具有特定意义的实体，并标注出其位置以及类型
  - **<ALB>中国有嘻哈</ALB>** **<PER>VAVA</PER>** 个性说唱 **<SONG>不想长大</SONG>**
- 实体类型
  - 人名、地名、机构名、产品名、专有名词等
- 研究领域
  - 新闻、Social Network Service、Query、生物、金融等
- 现有系统
  - Stanford、哈工大、Jieba、百度云、阿里云、腾讯云
- 模型：序列标注问题

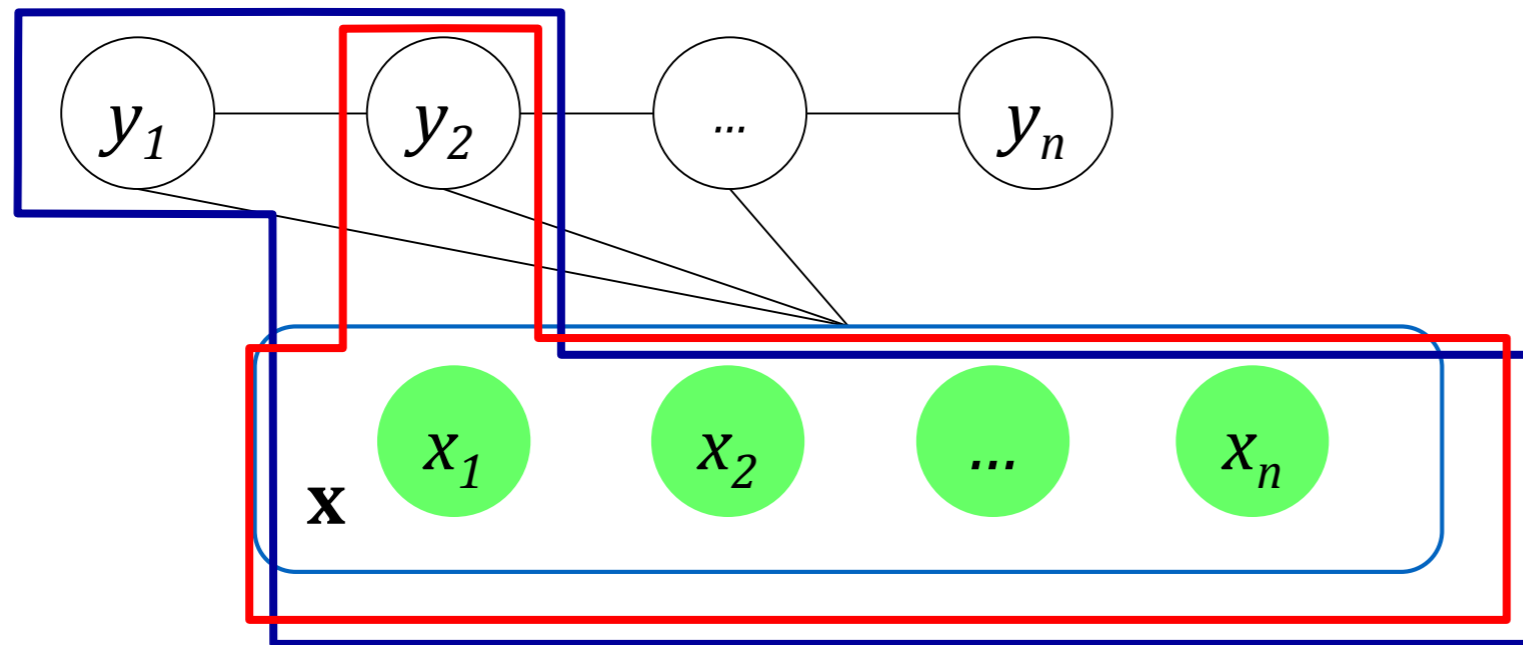
# 实体识别模型



# CRF

- 基于词/字的序列标注 ( BIEO )

中国	有	嘻哈	VAVA	个性	说唱	不	想	长大
B-ALB	I-ALB	E-ALB	B-PER	O	O	B-SONG	I-SONG	E-SONG



- # Unigram
  - U00:%x[-2,0]
  - U01:%x[-1,0]
  - U02:%x[0,0]
  - U03:%x[1,0]
  - U04:%x[2,0]
  - U05:%x[-1,0]/%x[0,0]
  - U06:%x[0,0]/%x[1,0]
  - U07:%x[-1,0]/%x[1,0]

- # Bigram
  - B

$$p(\mathbf{y}|\mathbf{x}) \propto \prod_{t=1}^n \exp \left\{ \sum_{l=1}^L \mu_l g_l(y_t, y_{t-1}, \mathbf{x}) \right\} \cdot \exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, \mathbf{x}) \right\}$$

# 视频领域实体识别

- 实体类型
  - 影视剧名、游戏名、音乐名、人名、角色名
- 挑战
  - 歧义性多
    - 功夫、十二生肖、长城、非诚勿扰
  - 规律性弱
    - 西游记之孙悟空三打白骨精
  - 别名多
    - 琅琊榜之风起长林、琅琊榜2、风起长林, etc.
  - 缺少语料



# 解决方法 (1)

- 词典

- 实体词

- 实时挖掘全网影视资料库 ( hourly update )

- 新词 ( 失恋哥、蓝瘦香菇 )

- 语料

- 通用和视频领域的准确人工标注语料
- 半监督的自动标注语料学习 ( Liu, et al. ACL2011 )
- 启发式规则

长城鹿晗成最怂士兵被马特达蒙踢飞

---

## Algorithm 1 NER for Tweets.

---

**Require:** Tweet stream  $i$ ; output stream  $o$ .

**Require:** Training tweets  $ts$ ; gazetteers  $ga$ .

```
1: Initialize  $l_s$ , the CRF labeler:  $l_s = train_s(ts)$ .
2: Initialize  $l_k$ , the KNN classifier:  $l_k = train_k(ts)$ .
3: Initialize  $n$ , the # of new training tweets:  $n = 0$ .
4: while Pop a tweet  $t$  from  $i$  and  $t \neq null$  do
5:   for Each word  $w \in t$  do
6:     Get the feature vector  $\vec{w}$ :  $\vec{w} = repr_w(w, t)$ .
7:     Classify  $\vec{w}$  with  $knn$ :  $(c, cf) = knn(l_k, \vec{w})$ .
8:     if  $cf > \tau$  then
9:       Pre-label:  $t = update(t, w, c)$ .
10:    end if
11:  end for
12:  Get the feature vector  $\vec{t}$ :  $\vec{t} = repr_t(t, ga)$ .
13:  Label  $\vec{t}$  with  $crf$ :  $(t, cf) = crf(l_s, \vec{t})$ .
14:  Put labeled result  $(t, cf)$  into  $o$ .
15:  if  $cf > \gamma$  then
16:    Add labeled result  $t$  to  $ts$ ,  $n = n + 1$ .
17:  end if
18:  if  $n > N$  then
19:    Retrain  $l_s$ :  $l_s = train_s(ts)$ .
20:    Retrain  $l_k$ :  $l_k = train_k(ts)$ .
21:     $n = 0$ .
22:  end if
23: end while
24: return  $o$ .
```

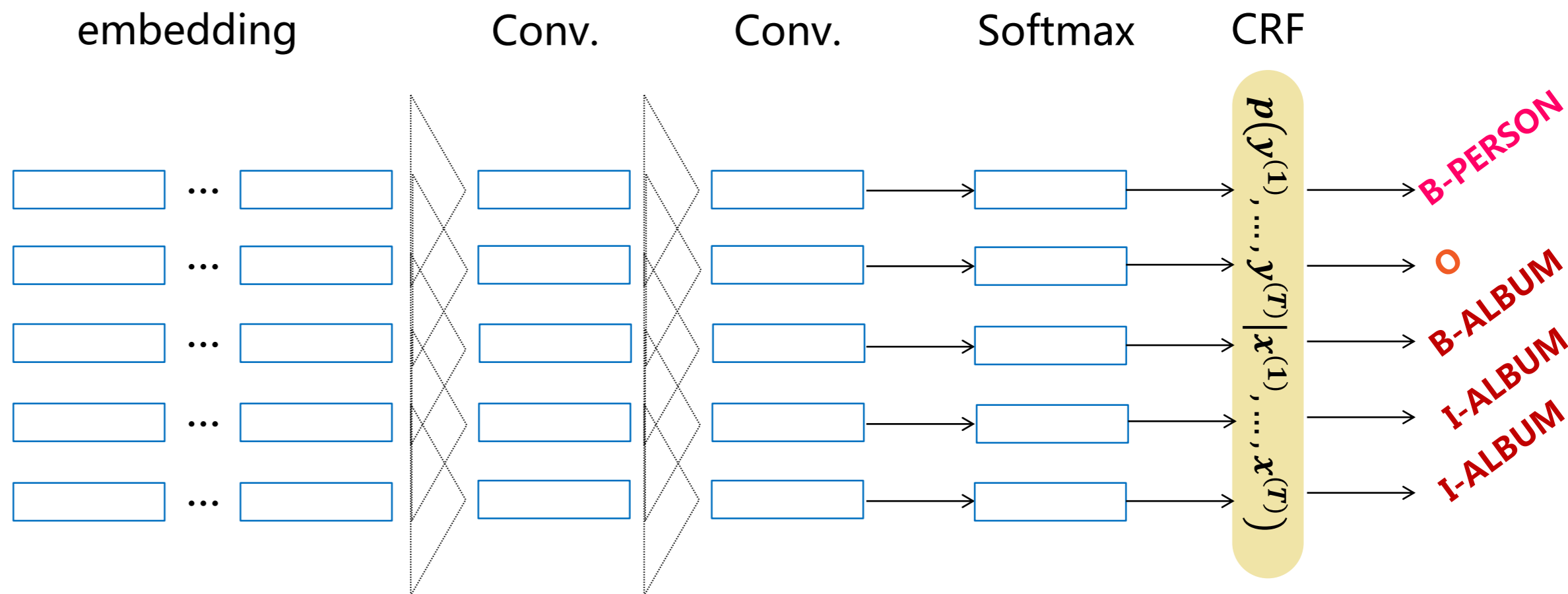


# 解决方法 (2)

- 特征与模型

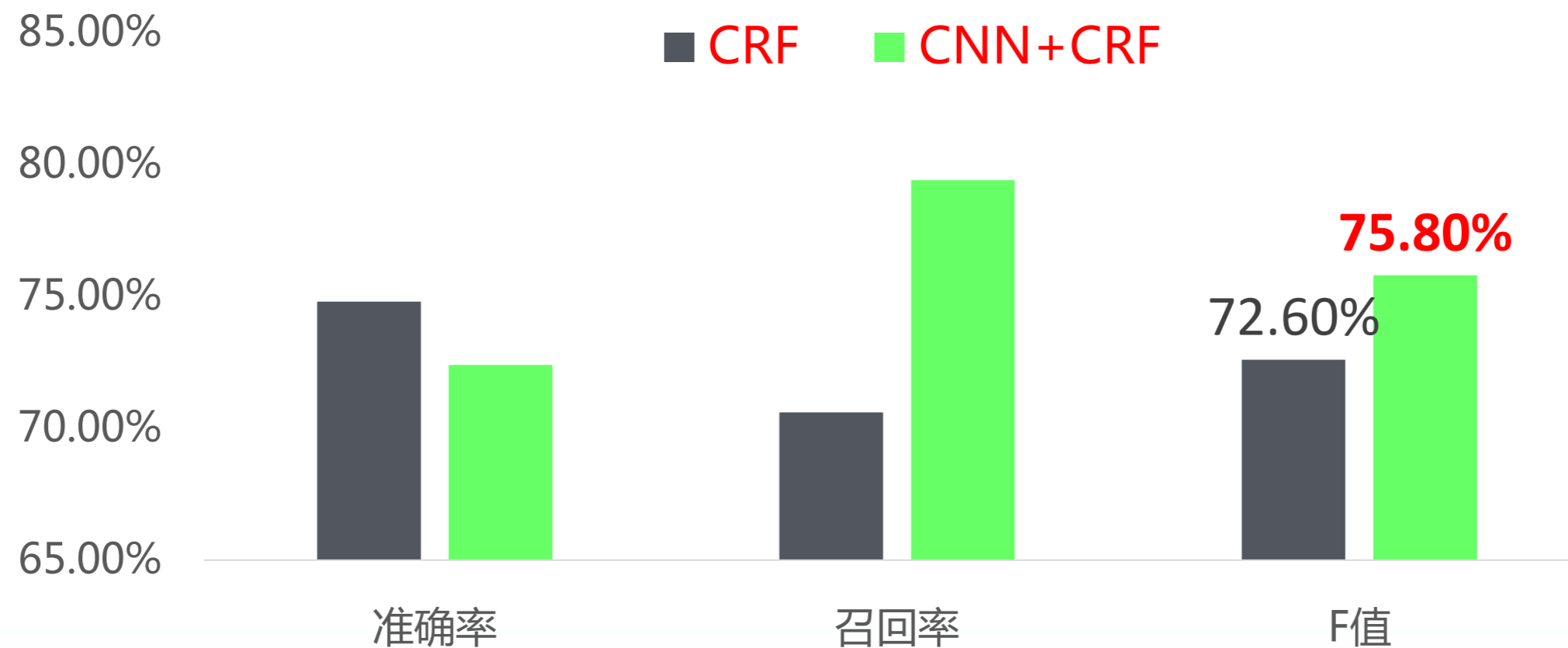
Word	POS	...	Brown Cluster
池子	$x_1^1(n)$	...	$x_1^L$
diss	$x_2^1(nx)$	...	$x_2^L$
中国	$x_3^1(ns)$	...	$x_3^L$
有	$x_4^1(v)$	...	$x_4^L$
嘻哈	$x_5^1(n)$	...	$x_5^L$

#L特征个数



# 实体识别性能评测

- 分词效果
  - 视频：91.21%、 微博：94.35%
- 影视剧名识别效果：CRF vs. CNN+CRF



# 标签

## Learning Deep Structured Semantic Models for Web Search using Clickthrough Data

Po-Sen Huang  
University of Illinois at Urbana-Champaign  
405 N Mathews Ave. Urbana, IL 61801 USA  
huang146@illinois.edu

Xiaodong He, Jianfeng Gao, Li Deng,  
Alex Acero, Larry Heck  
Microsoft Research, Redmond, WA 98052 USA  
{xiaohe, jfgao, deng, alexac, lheck}@microsoft.com

### ABSTRACT

Latent semantic models, such as LSA, intend to map a query to its relevant documents at the semantic level where keyword-based matching often fails. In this study we strive to develop a series of new latent semantic models with a deep structure that project queries and documents into a common low-dimensional space where the relevance of a document given a query is readily computed as the distance between them. The proposed deep structured semantic models are discriminatively trained by maximizing the conditional likelihood of the clicked documents given a query using the clickthrough data. To make our models applicable to large-scale Web search applications, we also use a technique called word hashing, which is shown to effectively scale up our semantic models to handle large vocabularies which are common in such tasks. The new models are evaluated on a Web document ranking task using a real-world data set. Results show that our best model significantly outperforms other latent semantic models, which were considered state-of-the-art in the performance prior to the work presented in this paper.

### Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval; I.2.6 [Artificial Intelligence]: Learning

### General Terms

Algorithms, Experimentation

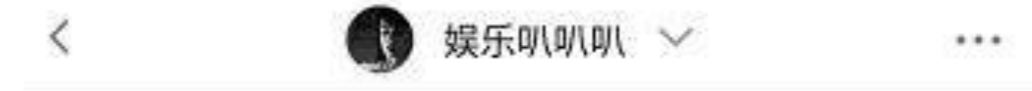
### Keywords

Deep Learning, Semantic Model, Clickthrough Data, Web Search

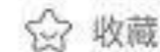
(LSA) are able to map a query to its relevant documents at the semantic level where lexical matching often fails (e.g., [6][15][2][8][21]). These latent semantic models address the language discrepancy between Web documents and search queries by grouping different terms that occur in a similar context into the same semantic cluster. Thus, a query and a document, represented as two vectors in the lower-dimensional semantic space, can still have a high similarity score even if they do not share any term. Extending from LSA, probabilistic topic models such as probabilistic LSA (PLSA) and Latent Dirichlet Allocation (LDA) have also been proposed for semantic matching [15][2]. However, these models are often trained in an unsupervised manner using an objective function that is only loosely coupled with the evaluation metric for the retrieval task. Thus the performance of these models on Web search tasks is not as good as originally expected.

Recently, two lines of research have been conducted to extend the aforementioned latent semantic models, which will be briefly reviewed below.

First, clickthrough data, which consists of a list of queries and their clicked documents, is exploited for semantic modeling so as to bridge the language discrepancy between search queries and Web documents [9][10]. For example, Gao et al. [10] propose the use of Bi-Lingual Topic Models (BLTMs) and linear Discriminative Projection Models (DPMs) for query-document matching at the semantic level. These models are trained on clickthrough data using objectives that tailor to the document ranking task. More specifically, BLTM is a generative model that requires that a query and its clicked documents not only share the same distribution over topics but also contain similar fractions of words assigned to each topic. In contrast, the DPM is learned using the S2Net algorithm [26] that follows the pairwise learning-



黄圣依也在第一时间发布微博表示自己会继续努力，也希望她能不忘当演员的初心，努力提升自己的演技拍摄更多更好的作品。



### 相关标签

黄圣依

演员的诞生

陶虹

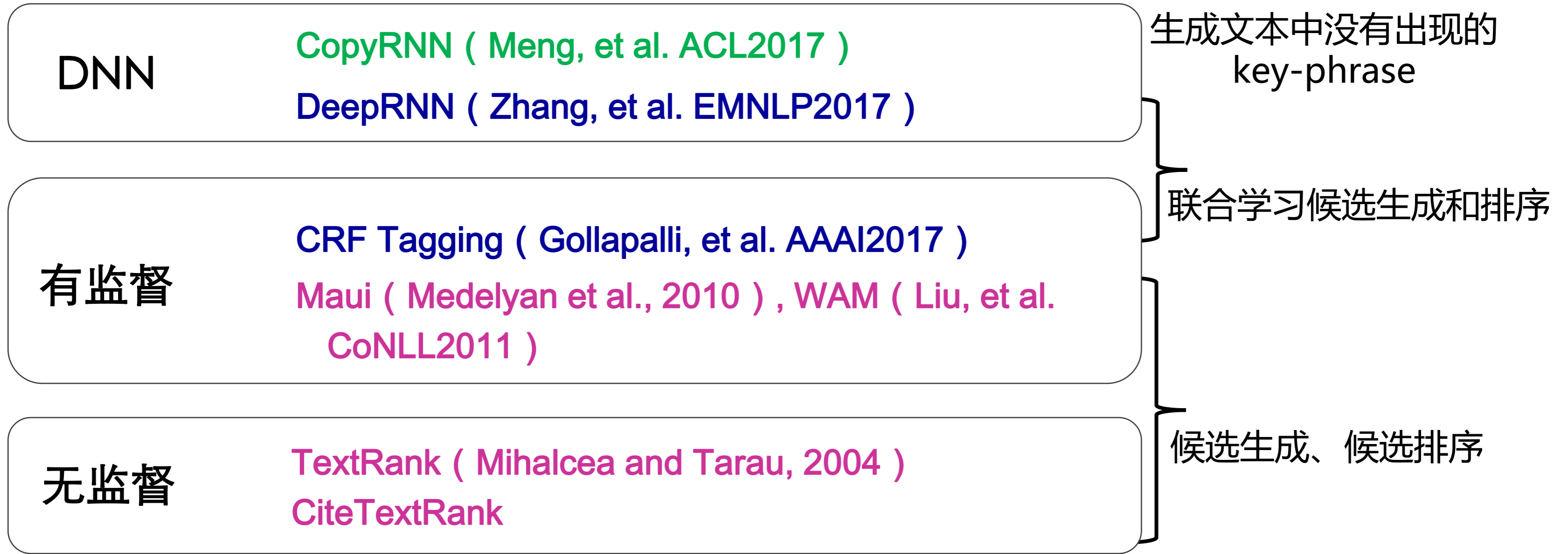
内地

刘烨

真人秀

新闻

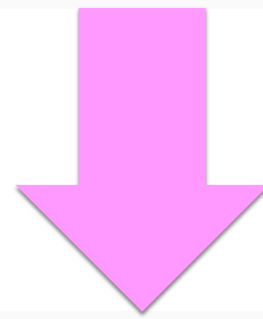
# 相关工作



# 相关工作

## VR影片《血肉与黄沙》斩获奥斯卡特别成就奖

.....第九届奥斯卡特别成就奖颁奖礼在美国洛杉矶举行，由著名导演亚利桑德罗·冈萨雷斯·伊纳里多（曾执导《鸟人》《荒野猎人》）拍摄的VR影片《血肉与黄沙》（Carne y Arena）获得了奥斯卡特别成就奖，成为首部获得奥斯卡奖的VR影片.....

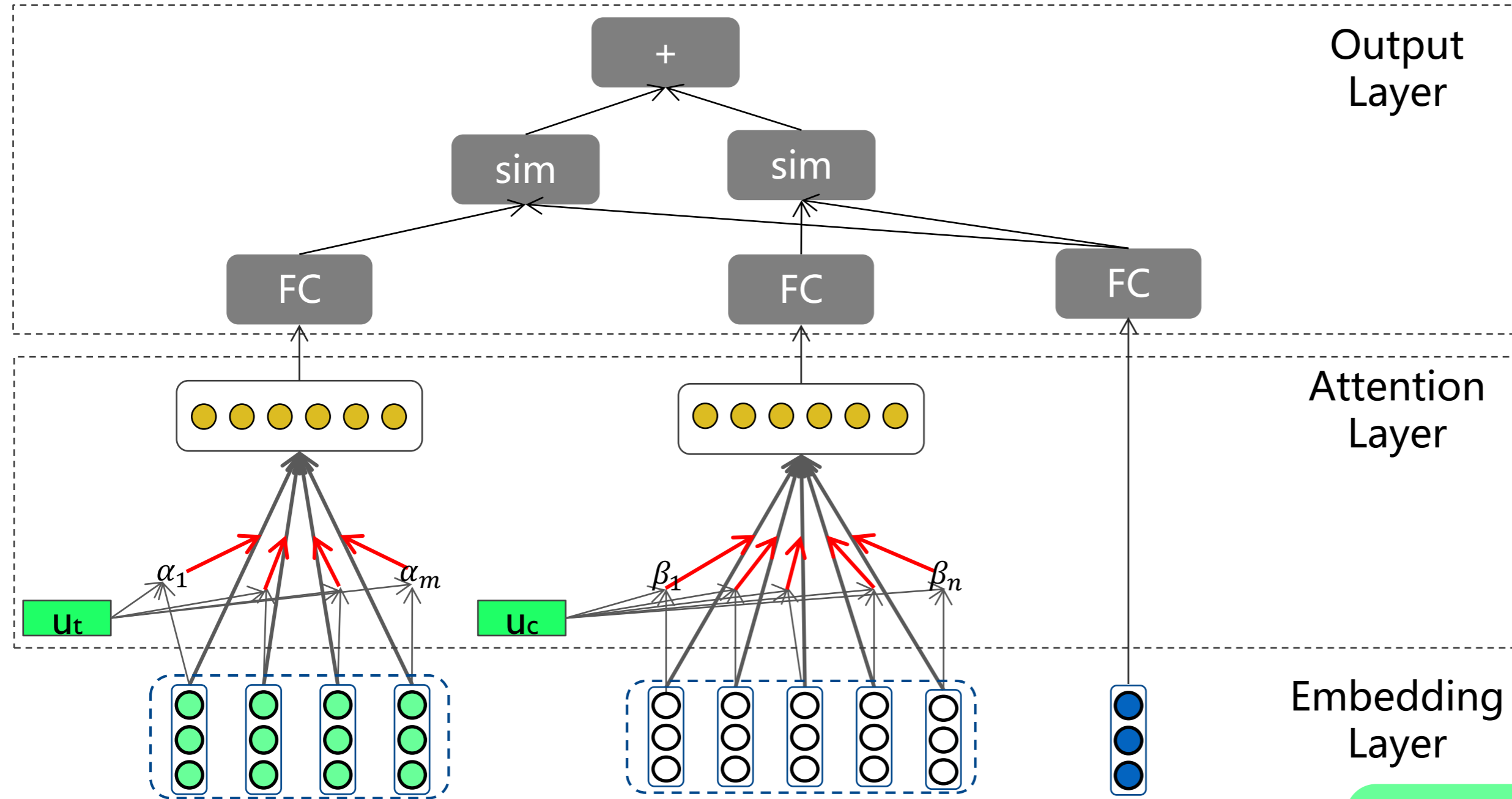


VR影片  
血肉与黄沙  
奥斯卡  
特别成就奖  
第九届奥斯卡  
颁奖礼

亚利桑德罗·冈萨雷斯·伊纳里多  
鸟人  
荒野猎人

VR 影片 《 血肉 与 黄沙 》 斩获 奥斯卡 特别 成就奖 ...  
B-KP I-KP O B-KP I-KP I-KP O O B-KP O O ...

# 注意力模型



- 可以抽取没有文本中出现的词组作为标签
- 更好地计算语义相关性

- 基于CRF的keyphrase抽取
- 标签库中的热门标签



VR影片《血肉与黄沙》斩获奥斯卡 近日，第九届奥斯卡特别成就奖颁奖典礼在美国洛杉矶...



# 标签性能

	召回率	准确率	F值
CRF tagging [Gollapalli, et al. 2017]	43.3%	62.3%	51.2%
CNN [Kim, et al. 2014]	52.3%	61.5%	56.5%
Our Attention Model	60.6%	70.5%	65.2%

# TABLE OF CONTENTES

- 理解视频内容
  - 中文词法分析
  - 预测（票房和流量）
- 理解视频用户
  - 舆情监测
  - 查询理解和意图搜索
- 总结



# 预测及应用场景

## 电视剧流量



版权采买，自制立项  
广告售卖，内容定级  
HCDN

- 长周期 ( **提前1年** )
- 中长周期 ( **半年** )
- 中短周期 ( **60日** )
- 播映中

## 电影票房



影业投资，版权采买  
内容宣发

- 长周期 ( **提前1年** )
- 中长周期 ( **半年** )
- 中短周期 ( **60、30日** )

动漫 综艺

...

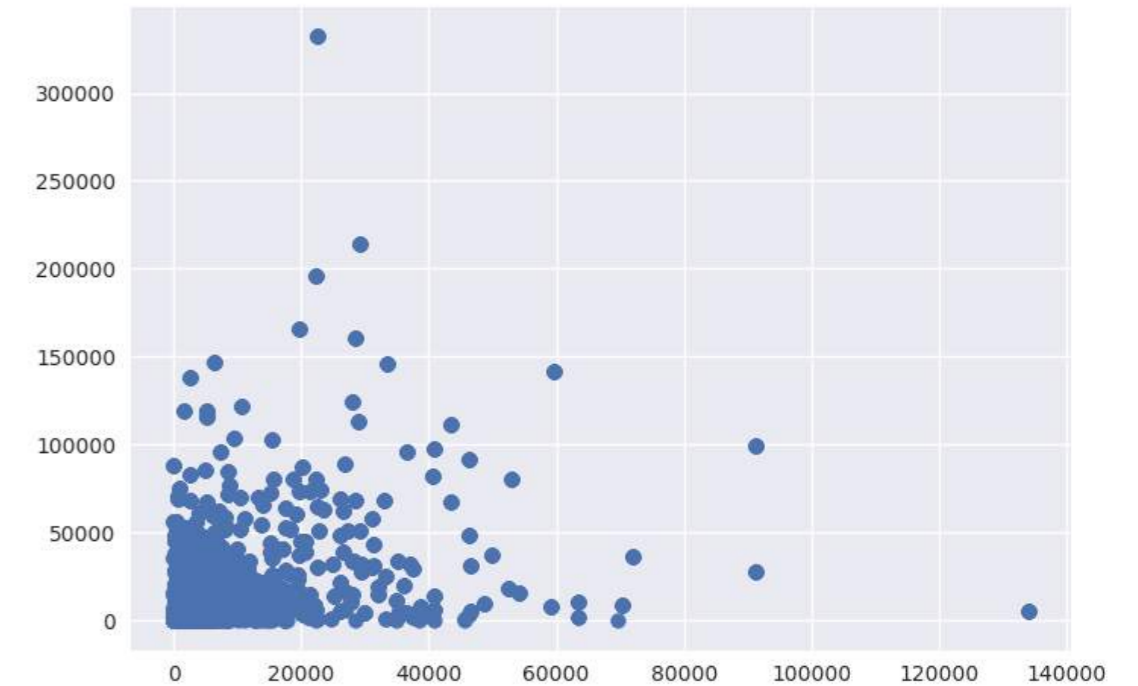
# 挑战与现状

- 挑战
  - 样本少
  - 影响因素多
    - 同档期竞争的影视剧、卡司突发事件、天气、电影排片率
  - 样片/剧本理解的难度
- 电影票房
  - 国外
    - Google ( 上映前一周 ) : 搜索、广告点击数据以及院线排片来预测票房 ( 2013 )
    - Epagogix ( 投资阶段 ) : 分镜头剧本中提取30,073,680个特定的评分指标
  - 国内
    - 猫眼 ( 上映影片 )、百度 ( 上映影片 )、艾漫 ( 投资 )

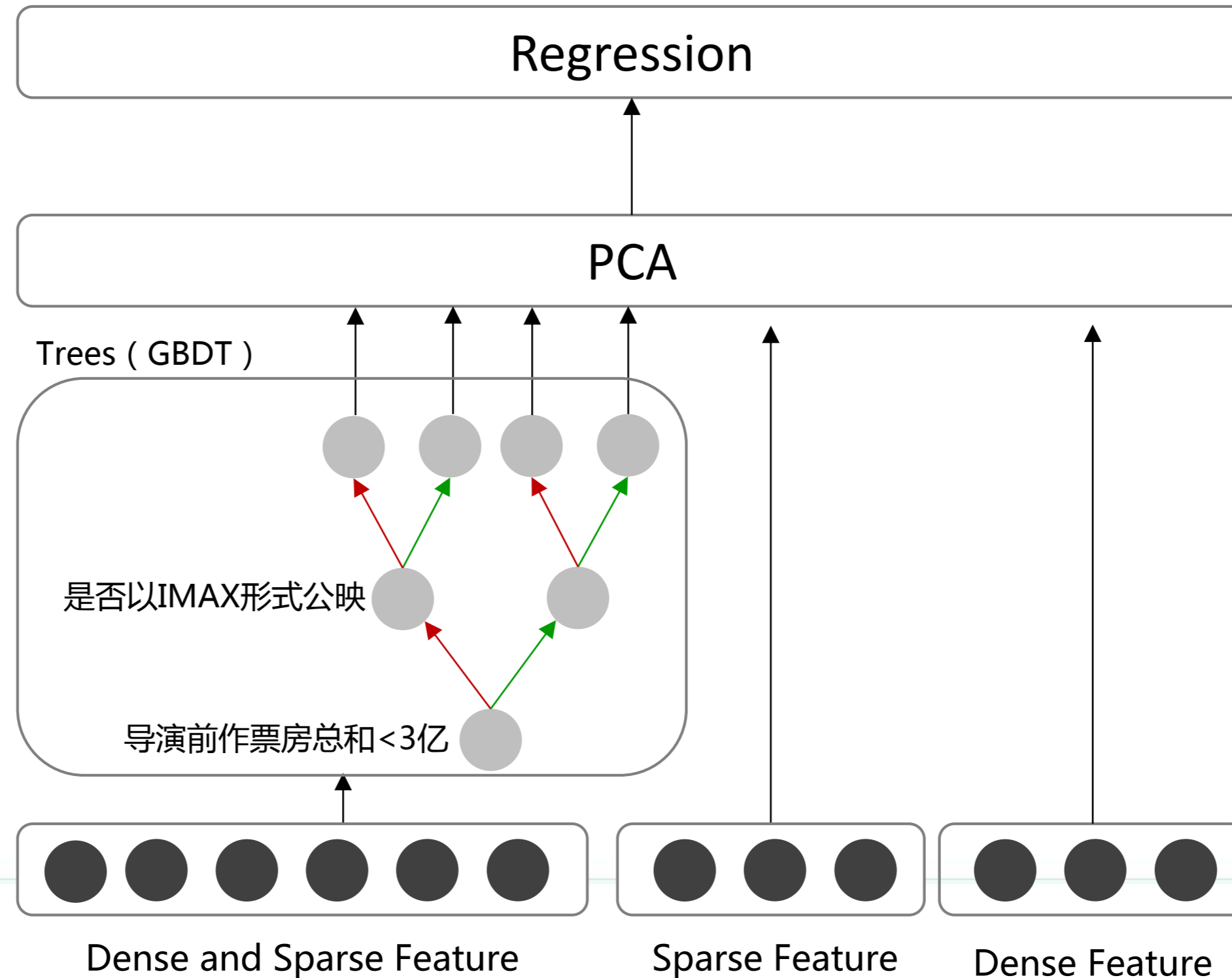
# 基于集成学习的预测模型

- 特征 ( **100+** )
  - 基本信息、历史前作、排播、搜索热度、社交热度
- 预处理
  - 缺失值处理、过采样、变换 ( log )
- 目标函数
  - 最小化log均方误差

$$loss = \frac{1}{2n} \sum_{i=1}^n (\log y_i - \log \hat{y}_i)^2 + \mu \|W\|_1$$

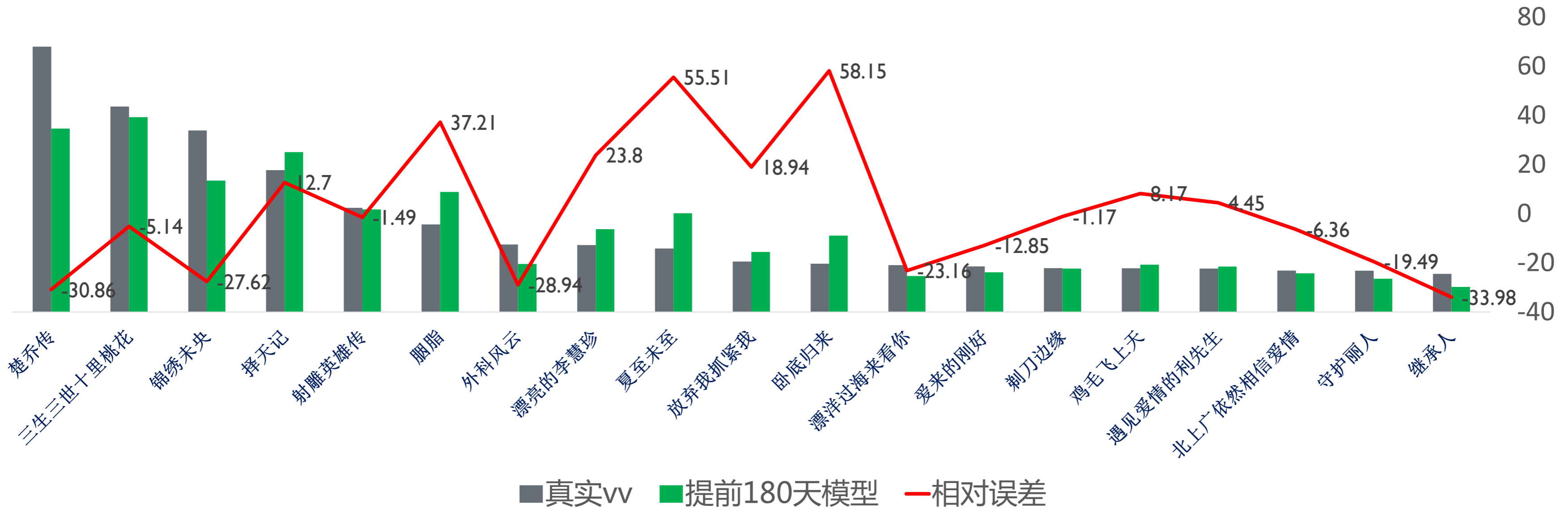


# 基于集成学习的预测模型



# 预测性能

- 电视剧流量预测R2准确率88%



# 预测性能

- 电影票房预测R2准确率81%



真实	4.34亿
预测	4.06亿
相对误差	6%

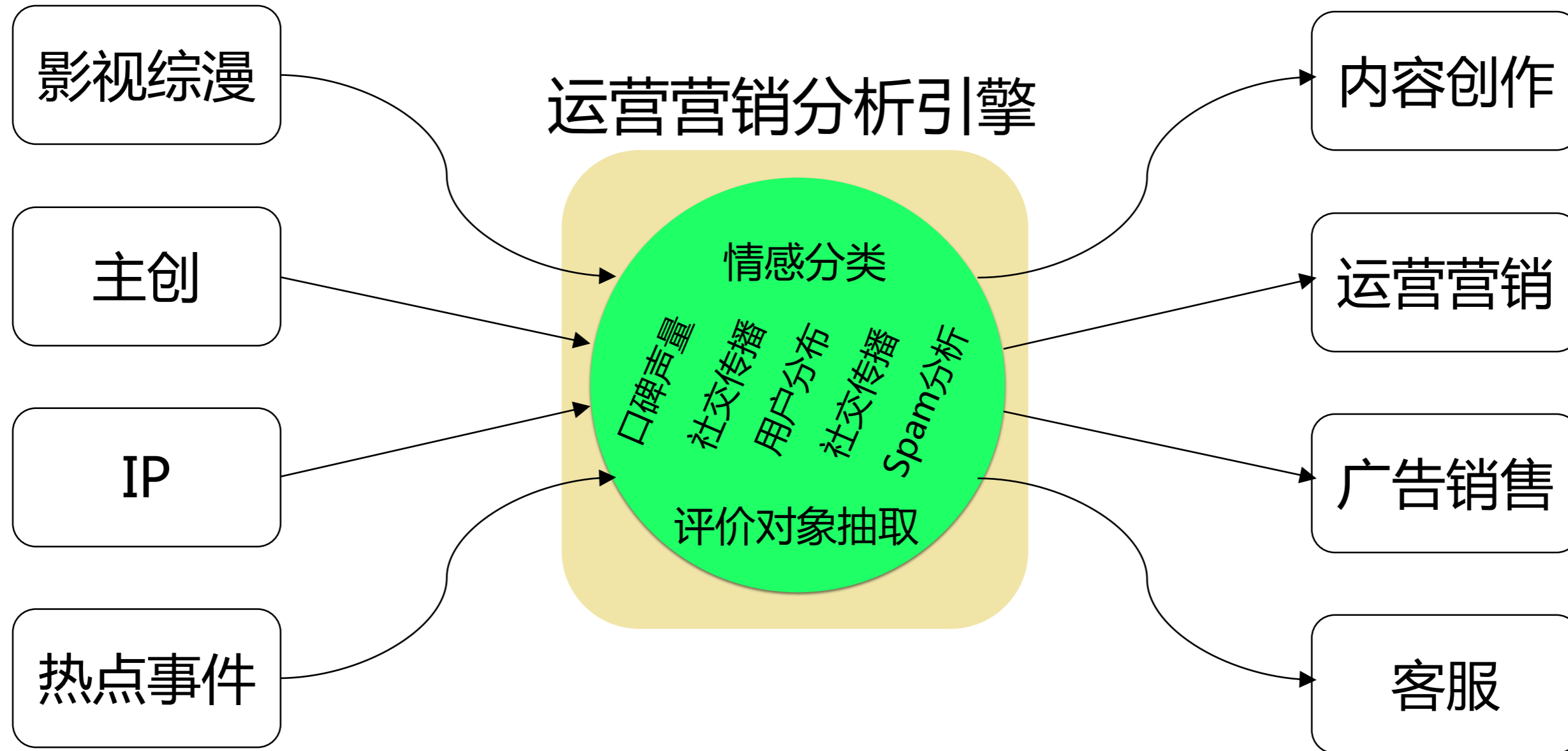


真实	13.27亿
预测	8.25亿
相对误差	37.8%

# TABLE OF CONTENTES

- 理解视频内容
  - 中文词法分析
  - 预测（票房和流量）
- 理解视频用户
  - 舆情监测
  - 查询理解和意图搜索
- 总结

# 运营营销





# 相关工作

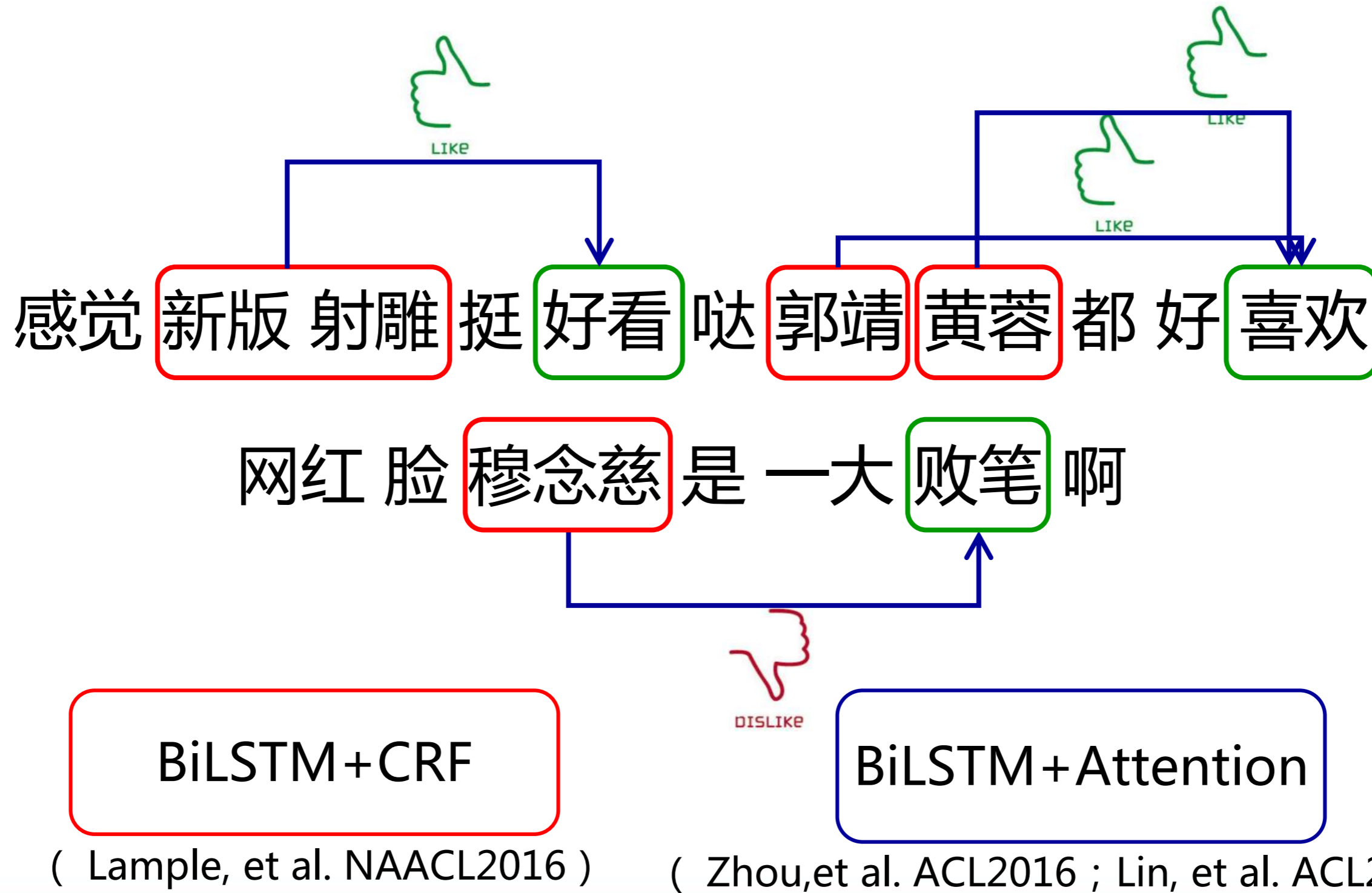
- COAE2016 ( 中文倾向性分析评测 )
  - Task A : 微博观点摘要
  - Task B : 影视评论的篇章级 - 句子级 - 词语级情感极性
- SemEval ( 2013~2017 ) : Sentiment Analysis in Twitter
  - Task A: Tweet Polarity Classification
  - Task B: Topic-based Tweet Polarity Classification
  - Task C: Tweet quantification

Tweet	Overall Sentiment	Topic-level Sentiment
Saturday without Leeds United is like Sunday dinner it doesnot feel normal at all (Ryan)	Weakly-Negative	Leeds United: Highly-Positive

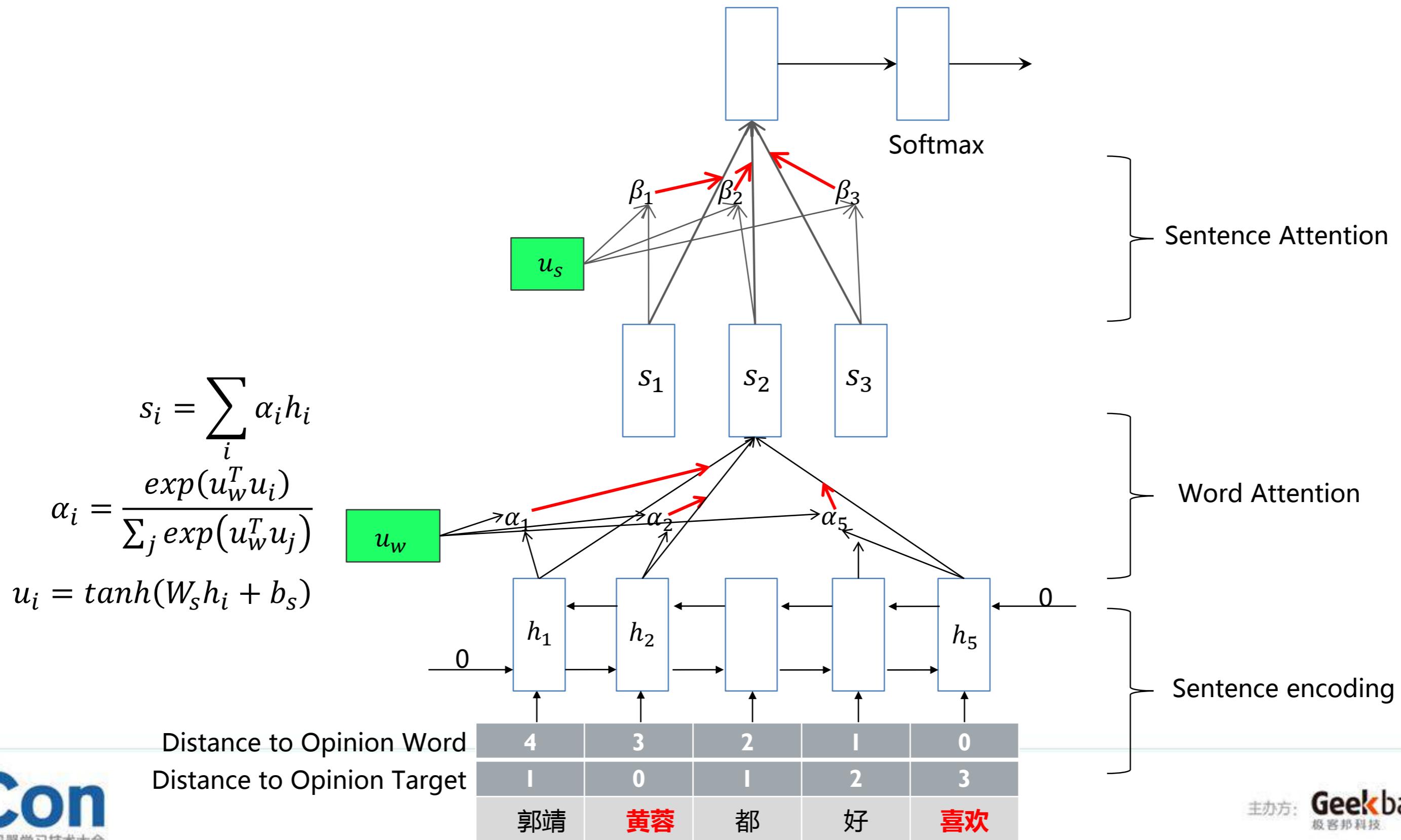
# 视频场景下舆情分析任务定义

- Given  $(u_i, a_i, c_i)$  :  $u_i$ =user,  $a_i$ =album,  $c_i$ =评论
  - Spam
    - 水军账号识别
    - Spam帖子识别
  - Paragraph/Sentence
    - 识别段落/句子的观点倾向性 : CNN、Bi-LSTM ( **83%** )
    - **挑战 : 郭靖 黄蓉 都好喜欢网红脸 穆念慈 是一大败笔 啊**
  - Phrase ( Aspect )
    - 评价对象、评价词抽取 ; 关系抽取 ; 评价对象聚合

# Phrase级舆情抽取



# BiLSTM + Hierarchical Attention



# DEMO



# DEMO



# TABLE OF CONTENTES

- 理解视频内容
  - 中文词法分析
  - 预测（票房和流量）
- 理解视频用户
  - 舆情监测
  - 查询理解与意图搜索
- 总结

# 查询理解





# 意图搜索

王菲的女儿

胡歌古装电视剧大全

演过黄蓉的演员有哪些

邓超老婆演过的电视剧

全网 王菲的女儿 取消

相关 最新 最热 高级筛选

王菲的女儿

王菲 女儿 窦靖童 李嫣

说明: 王菲的女儿是窦靖童和李嫣

王菲女儿窦靖童被骂长得丑 李嫣维护姐姐发文呛网友  
发布时间: 2015-10-13  
上传者: 1613581573

王菲放弃裂唇女儿被痛骂, 可李嫣说的道出了实情  
发布时间: 2017-07-19  
上传者: raul16\_1802

王菲女儿  
共27个视频  
简介: 在王菲上海站的演唱会上, 小女儿李

全网 胡歌古装电视剧大全 取消

相关 最新 最热 高级筛选

琅琊榜 风中奇缘 风中奇缘DVD版

仙剑奇侠传3 神话 仙剑奇侠传1

查看全部

全网 演过黄蓉的演员有哪些 取消

相关 最新 最热 高级筛选

演过黄蓉的演员有哪些

杨明娜 李一桐 魏秋桦 翁美玲  
孔琳 林依晨 朱茵 宋宁

胥渡吧配音演员现场给《射雕英雄传》黄蓉配音, 真的太像了  
发布时间: 2017-10-01  
上传者: 影视配音湿

演员拍鬼片, 谁知混进来一个真鬼, 拍出了逼真的现场  
发布时间: 2017-04-28  
上传者: a7man1314\_1703

《射雕英雄传》黄蓉的配音演员

全网 邓超老婆演过的影视剧 取消

相关 最新 最热 高级筛选

邓超老婆演过的影视剧

恶棍天使 大染坊 甜蜜蜜

小姨多鹤 分手大师 关云长

查看全部

# TABLE OF CONTENTES

- 理解视频内容
  - 中文词法分析
  - 预测（票房和流量）
- 理解视频用户
  - 意图搜索
  - 舆情监测
- 总结

# 总结

- 理解视频内容
  - 中文词法分析
  - 智能标签
  - 票房和流量预测
  - 热点事件发现和摘要
  - 智能审核（审核、标题党、软色情等）
- 理解视频用户
  - 舆情监测
  - Query Understanding
  - 助手与客服

Thanks!