

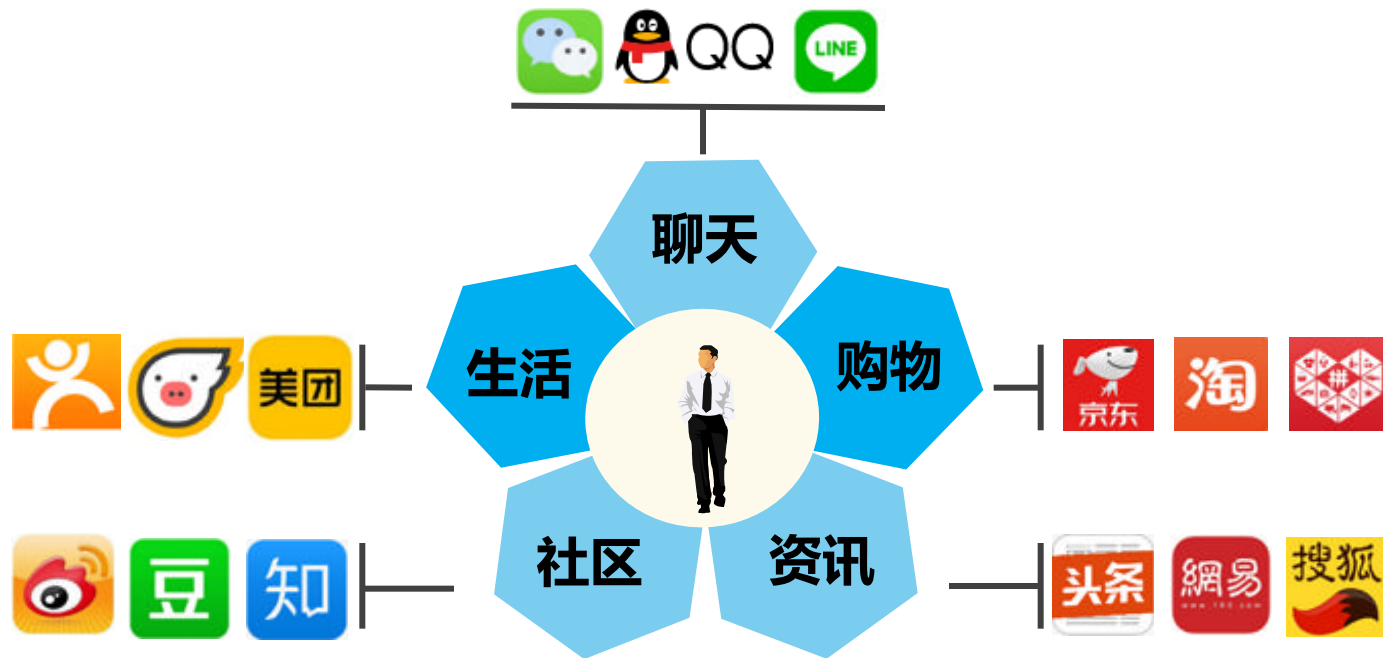
# 社交媒体情感分析



中国科学院 信息工程研究所  
INSTITUTE OF INFORMATION ENGINEERING, CAS

# 社交媒体

社交媒体是人们用来创作、分享、交流意见、观点及经验的网络平台。社交媒体已经涉及到现代人生活的方方面面，成为信息传播和维系社会关系的重要渠道。而文本是社交媒体交流的主要载体。



# 社交媒体中的情感

## 社交媒体上存在大量包含用户情感的文本

七月的色烈芬\_ 看过 ★★★★★

19603 有用

开篇长镜头惊险大气引人入胜 结合了水平不俗的快剪下实打实的真刀真枪 让人不禁热血沸腾 特别弹簧床架挡炸弹 空手接碎玻璃 弹匣割喉等帅得飞起！就算前半段铺垫节奏散漫主角光环开太大等也不怕 作为一个中国人 两个小时弥漫着中国强大得不可侵犯的氛围 还是让那颗民族自尊心砰砰砰跳个不停。

人民日报

人民日报 V

【最后倒计时，我们准备好了！#受阅官兵集结完毕#等待检阅】北京，天安门，此刻，世界为之瞩目！#国庆大阅兵#受阅部队整装待发，铁甲铮铮！今天，无论你在何处，请转发微博，为阅兵喝彩，为中国点赞！



宝宝\_小仙女 Lv6 VIP

口味：不错 环境：满意 服务：满意 人均：160元

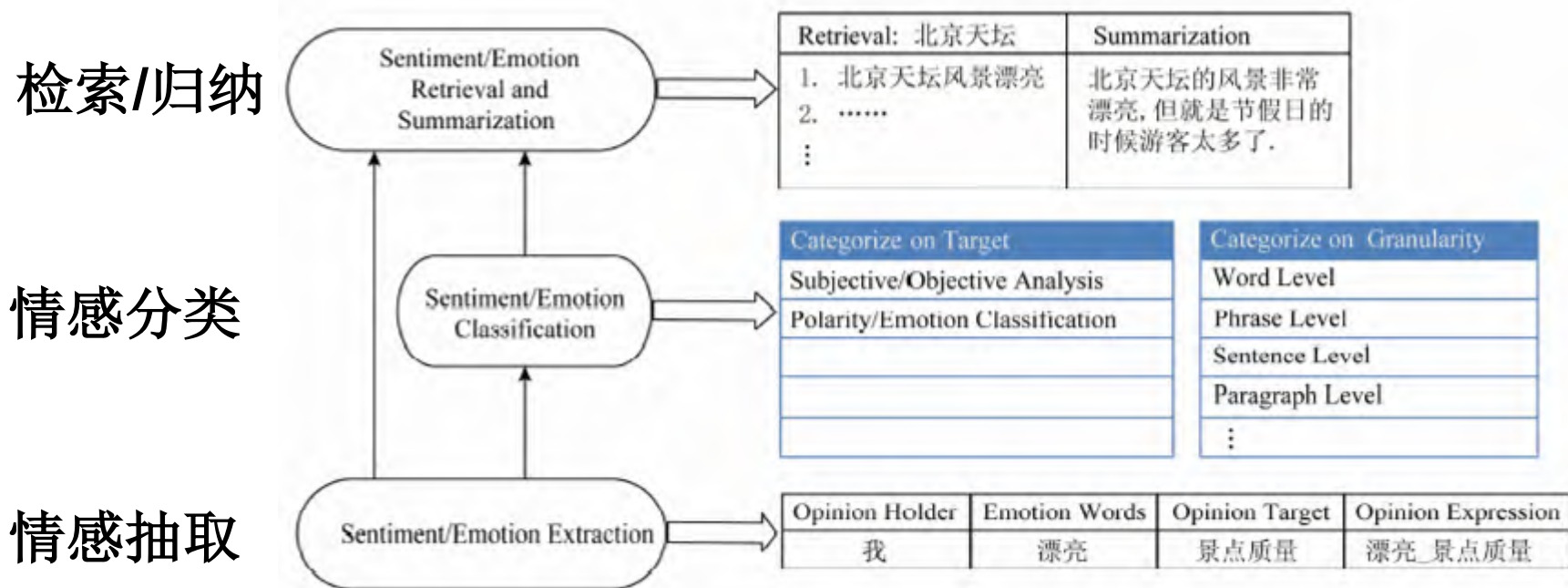
位置：中华老字号，这家店位于前门商业街，人气十分旺

环境：招牌富丽堂皇，内堂也很华丽，金光灿灿，古色古香，环境干净卫生。

服务：服务员十分热情好客，服务到.....

# 社交媒体中的情感分析

情感分析是一种重要的信息组织方式，研究的目的是自动挖掘和分析文本中的**立场、观点、看法、情绪和喜恶**等主观信息。其一般的研究框架包含**情感抽取、分类、检索与归纳**等任务



情感分析研究框架

# 情感分析具有巨大的应用价值

当前，随着微信、微博、论坛和社交网络等应用的兴起，社交网络上汇聚着海量的信息。情感分析在**社会管理、商业决策、信息预测**等各个反面有着广泛而重要的应用价值



**社会管理**

通过情感分析进行舆情态势感知：

- 发现民众意见倾向，客观反映舆情状态
- 发现非理性情绪的“群体极化现象”

# 情感分析具有广泛的应用场景



商业决策

通过情感分析处理评论信息：

- 商家改进产品
- 消费者确定购买意向



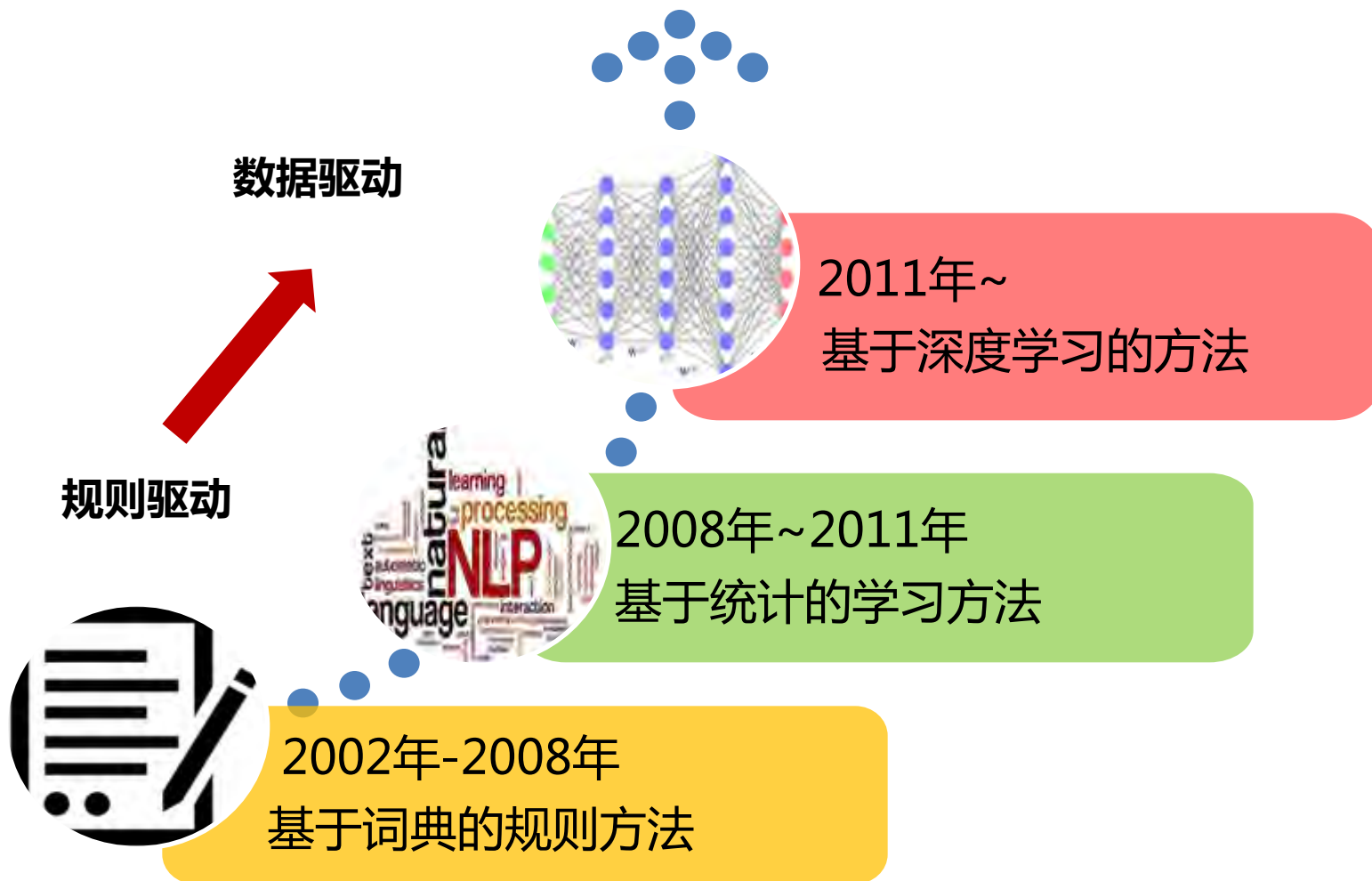
信息预测



通过情感分析：

- 商业预测，Twitter可以一定程度上预测3-4天的股市变化
- 选情预测，制定宣传策略

# 情感分析研究历史



# 基于词典的规则方法

- 早期的情感分析工作主要围绕主客观判定和情感极性判定两方面展开
- 2003年有学者阐述了一种情绪分析方法，用于从文档中提取与**特定主体的正面或负面相关的情绪**，该方法通过使用语法分析器和**情感词典**进行语义分析

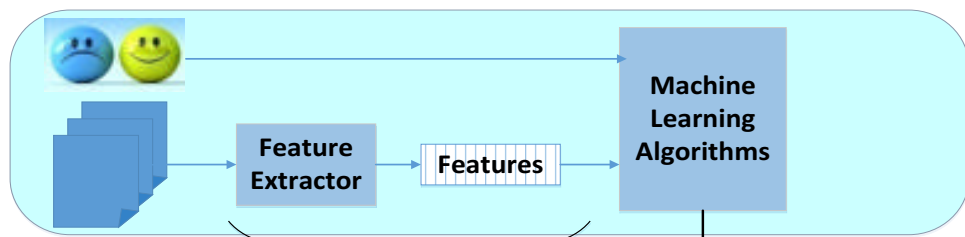




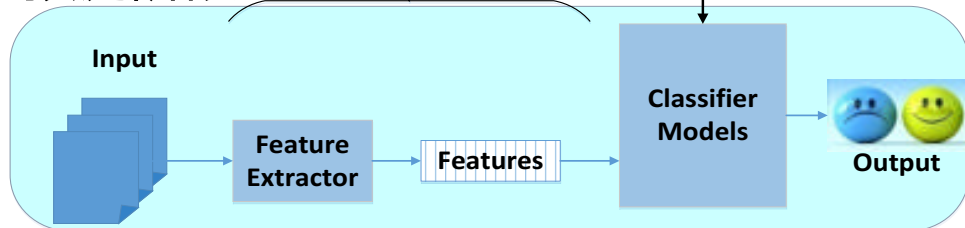
# 基于统计机器学习的方法

机器学习方法主要是采用**有监督的学习方式**，在有标注的训练语料上训练一个情感分类器，然后用于未标注数据的情感极性预测

训练阶段



预测阶段



- One-hot vector
- N-grams
- Lexicons
- Patterns
- ...

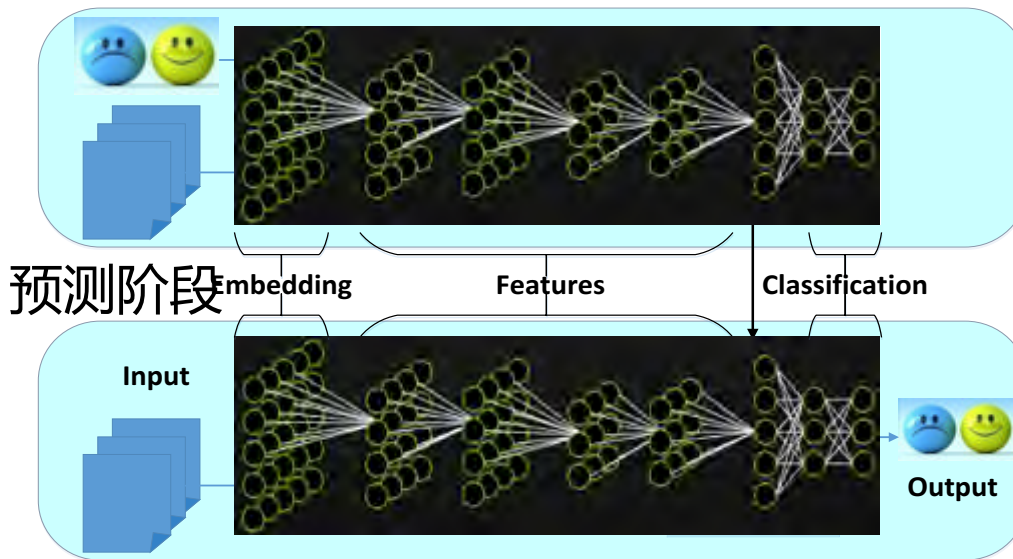
缺点：

- 过于依赖标注数据
- 人工提取特征

# 基于深度学习的方法

基于深度学习的方法文本特征是**自动提取**的，通过神经网络学习文本中所蕴含的语义信息，达到情感分析的目的

训练阶段

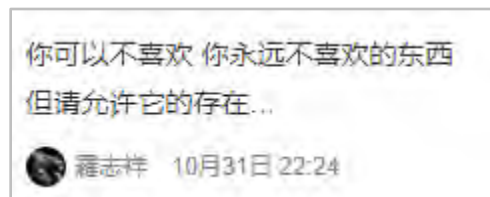


- 利用大量无标注数据进行预训练，减少对标注数据的需求
- 解决传统方法的稀疏性问题
- 自动学习特征表示，更灵活

# 社交媒体情感分析面临的挑战

## 数据稀疏问题

### — 社交媒体文本数据特征稀疏



### — 社交媒体中新词层出不穷



盘他



好嗨哦!



咱也不敢说  
咱也不敢问



我太难了

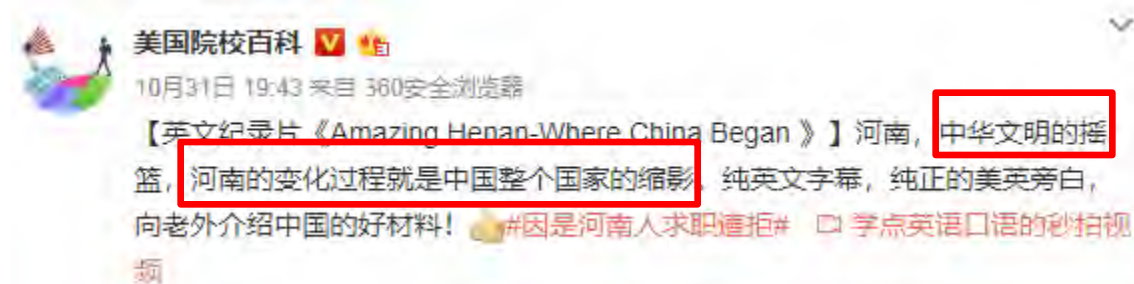
### — 社交文本存在很多噪音

样本	输入	输出
原始样本	never coming here again found hair in my to go food.	<b>Negative</b>
噪音样本	enver coming here again found hair in my to go food.	<b>Positive</b>

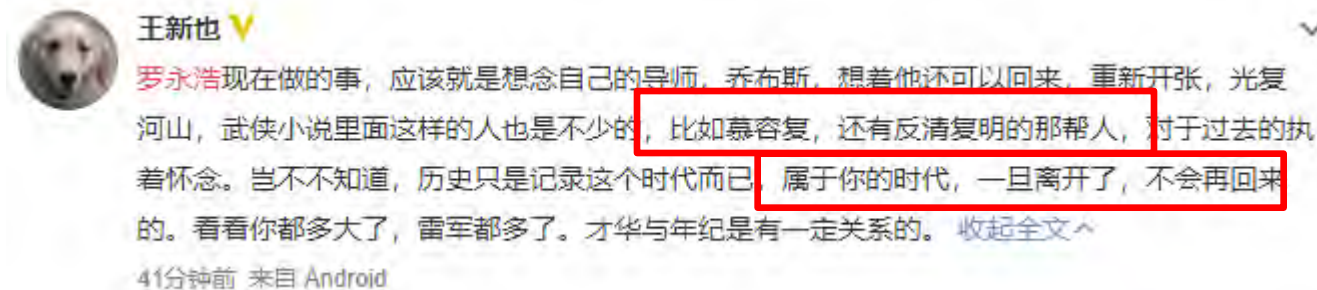
# 社交媒体情感分析面临的挑战

## • 隐式情感推断困难问题

- 社交媒体中文本情感表达方式复杂，多数没有显式情感词



- 社交媒体上存在很多讽刺与挖苦的语言现象，这种隐式情感（言外之意）的推断会比显式情感分析（文字表面含义）更加困难



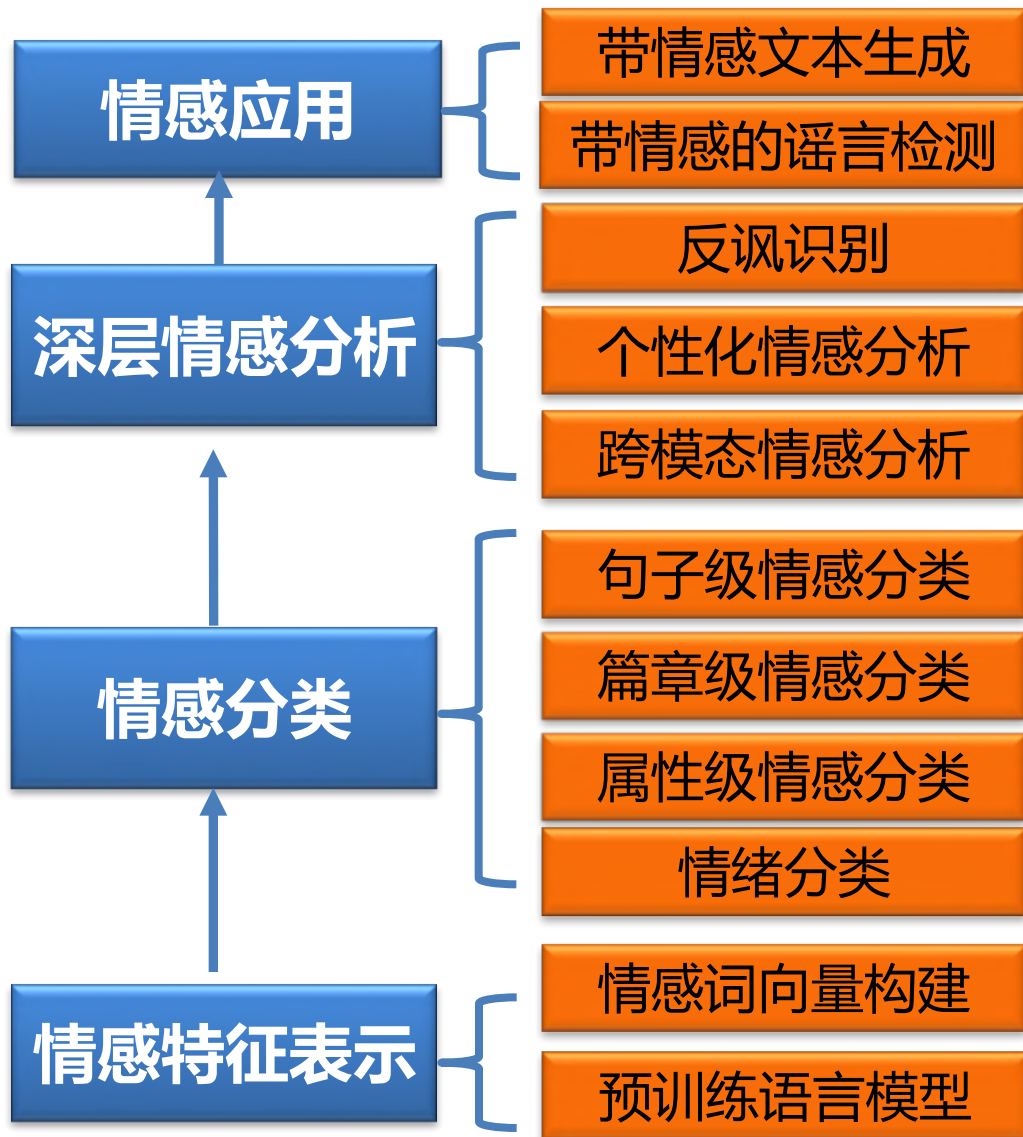
# 社交媒体情感分析面临的挑战

- 社交媒体数据规模大、数据跨模态问题

互联网上每天都在生成大量的情感文本和音视频，对海量数据进行挖掘需要解决数据特征规模大、特征复杂度高、数据时效性高、数据跨模态、模型训练频繁等问题。



# 社交媒体情感分析关键技术



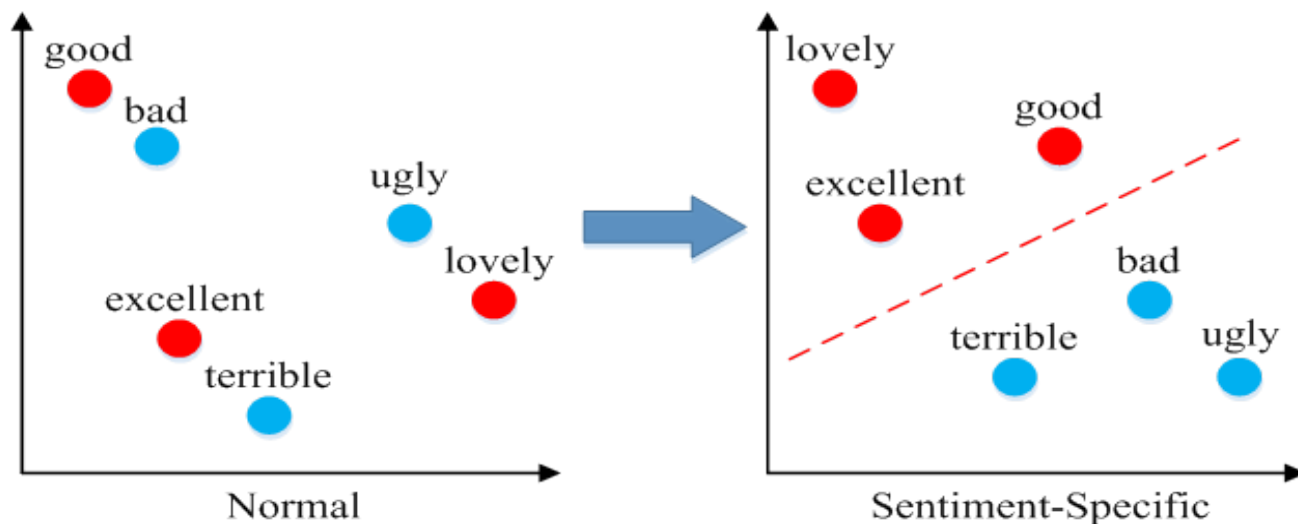
---

# 一、情感词向量构建

# 研究动机

传统词向量学习方法建立基于目标词上下文的窗口语言模型。这样的词向量包含了词的语义和语法信息，但也会将具有相似上下文但是情感极性相反的词映射到相近的向量空间，导致**很难直接用于情感分析**

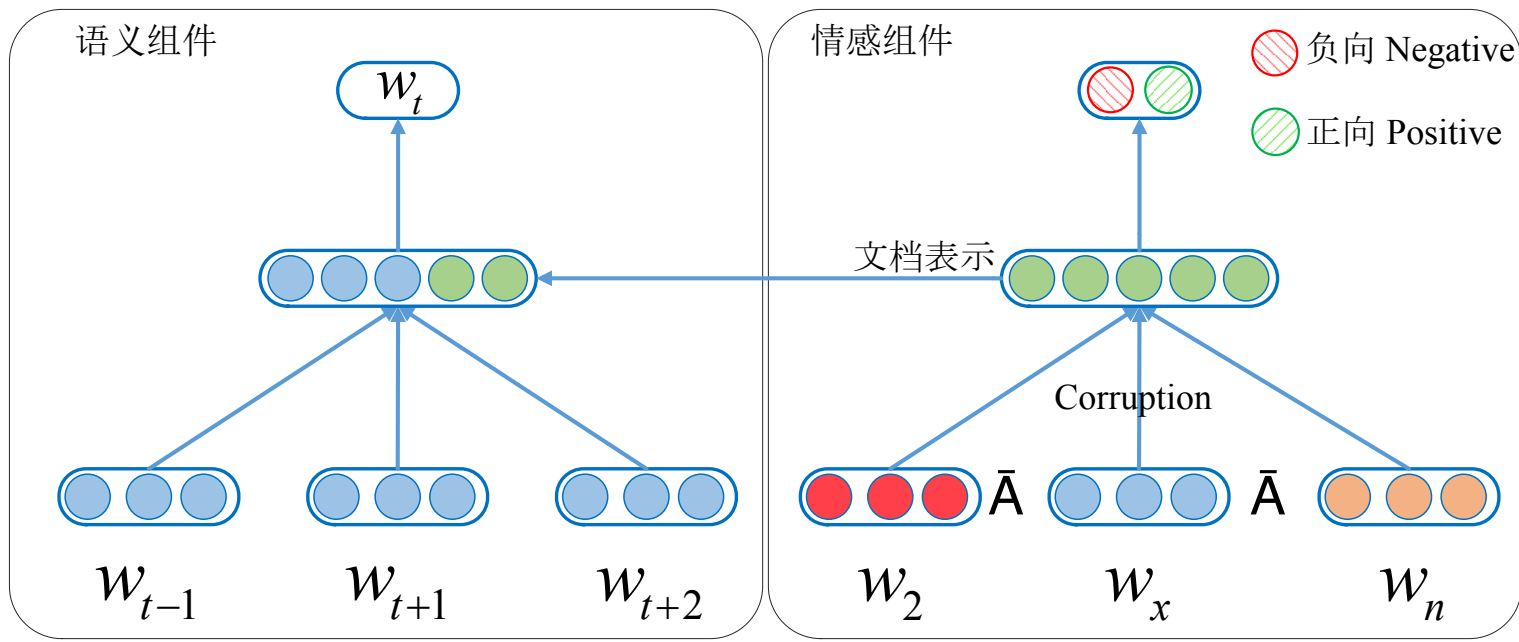
*I feel **good** for it.*  
*I feel **bad** for it.*





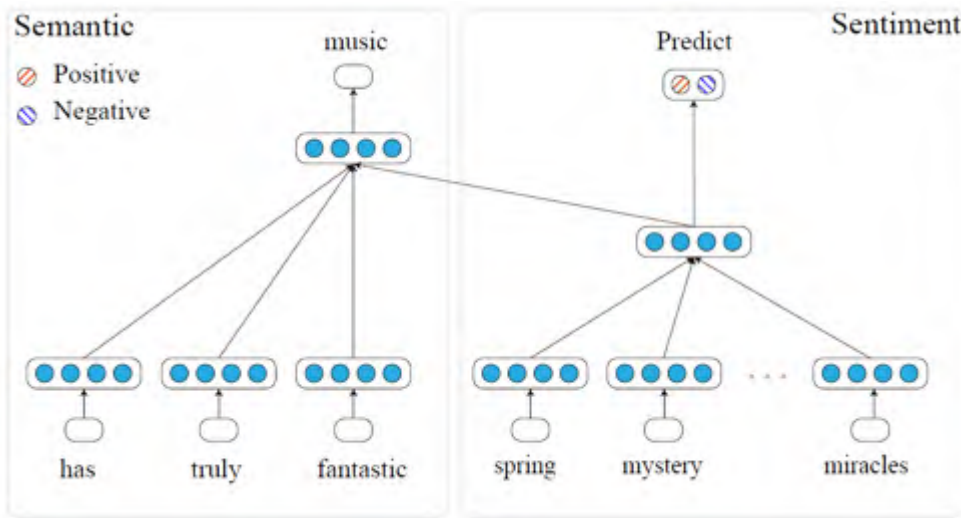
# 情感词向量学习方法

- 模型包括两个组件
  - **语义组件**：通过无监督的方式学习词的语法和语义信息
  - **情感组件**：通过有监督的方式，结合有偏的弃词机制，学习全局的文档情感语义表示
- 整个框架通过目标词的上下文和全局的文档情感倾向性，共同预测目标词，从而学习到具有情感属性的词向量

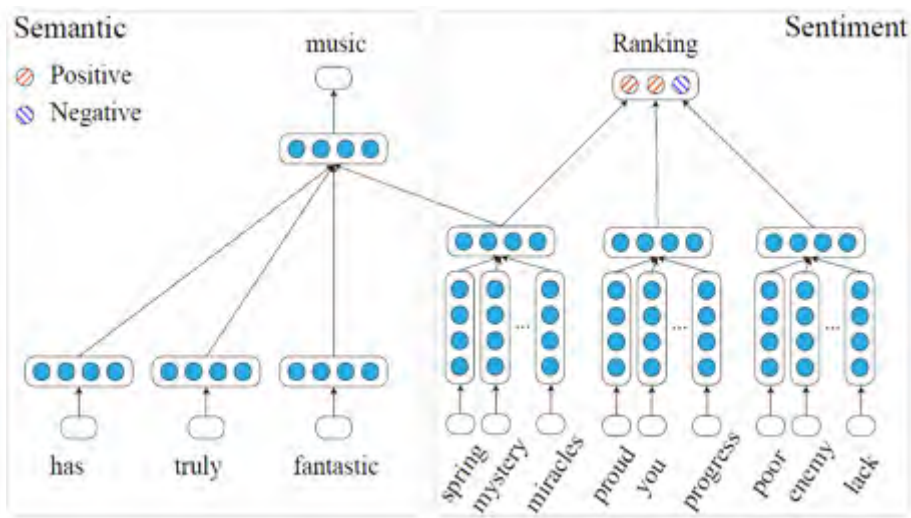


# 情感词向量学习模型

- **预测模型**：预测模型将全局文档的情感极性分布问题看作是分类问题，通过有监督学习对全局文本表示进行情感二分类
- **排序模型**：排序模型的基本思想是使情感极性一致的文档尽可能相近，不一致的文档尽可能远

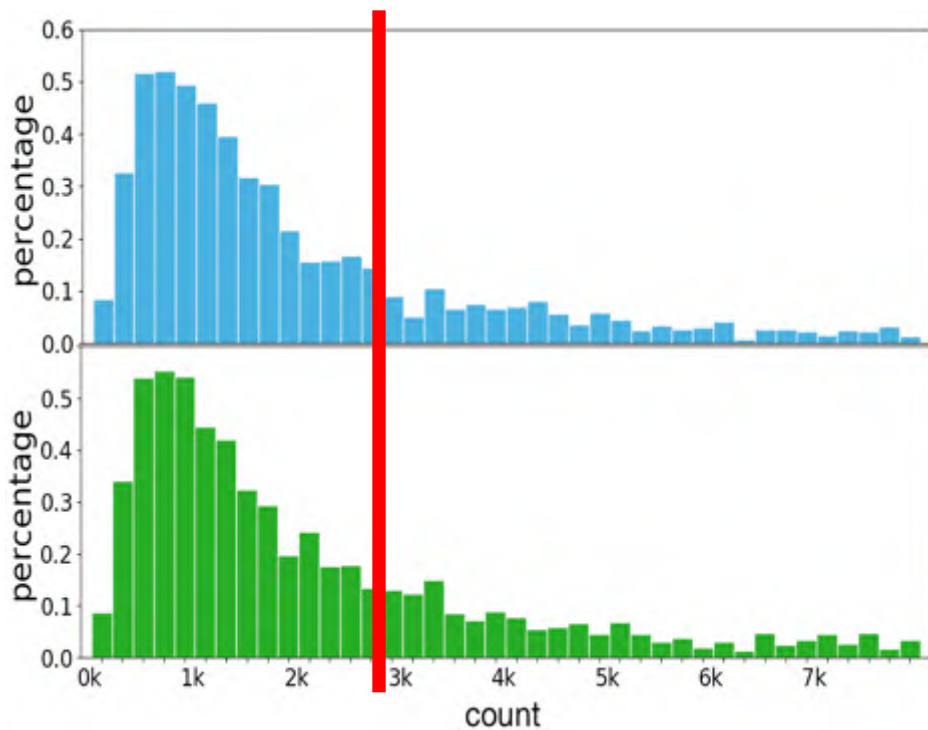


预测模型



排序模型

# 采样弃词机制



情感词在不同词频范围所占的比例  
IMDB (上) Twitter (下)

➔ 对语料中的情感词进行统计分析，发现大多数情况下，情感词的词频要低于普通词的词频

- ✓ 在构造文档表示时，可以根据词频信息有偏的丢弃一些词高频词，从而使得模型偏好于情感词比例高的词频范围
- ✓ 可以加速模型的训练，减小模型复杂度

---

## 二、情感分类

# 情感分类任务定义

**情感分类**：输入句子/篇章，输出句子/篇章的褒贬极性



每日电影头条

少年的你 ★★★★★，过分了！《少年的你》简直是国产电影的巅峰之作！135分钟每一秒每一分都完美到爆！一点都不觉得长，电影结束了坐在座位上都不想离去，导演把节奏把握的太好了！全程温暖哭泣！曾国祥快出来受夸，太完美了，出乎意外的拍摄手法，前所未见，周冬雨全程演的很棒！不愧是好演员！而易烱千玺！天呐！世界上怎么会有易烱千玺这样的男孩子



**情绪分类**：输入句子/篇章，输出句子/篇章的情绪类别（喜怒哀乐悲恐）



李易峰

很高兴昨晚和CaraDelevingne一起参加了@TAGHeuer泰格豪雅 #MONACO50#终极派对，晒一下新鲜出炉的时髦D.N.A.“验血报告”，跟我“同年代”的举手！



# 句子级情感分类

- 句子级情感分析是指对**单句**的情感极性进行分析判断的任务
- 基于深度神经网络的方法主要是对句子所包含的**语义信息**进行表示，进而对其情感极性进行判别
- 常用的方法有基于**CNN**、**RNN**、**Recursive-NN**的方法以及最新的**Transformer**等

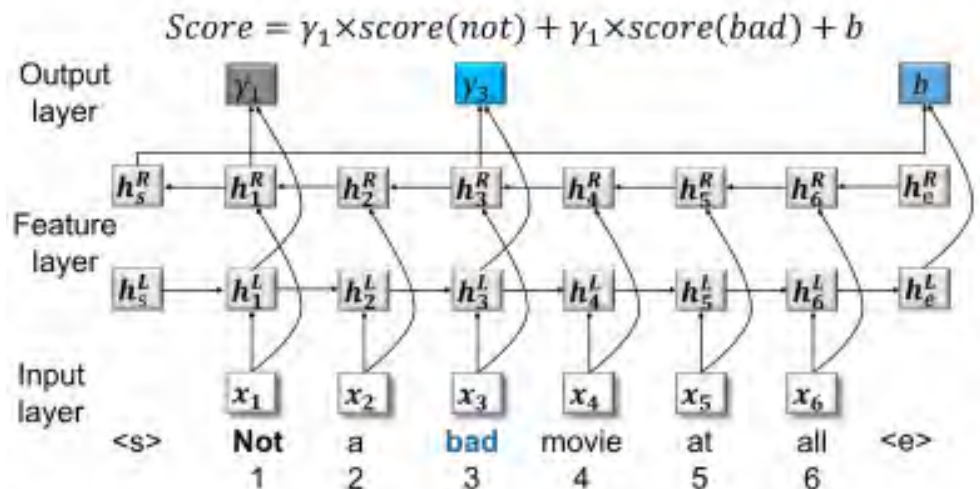
The performances were all **really** It's an **insignificant** **[criticism]**<sub>-1→-0.5</sub>.

**Nobody** gives a **[good]**<sub>+3→-1</sub> performance in this movie

She's **not** **[terrific]**<sub>+5→+1</sub> but **not** **[terrible]**<sub>-5→-1</sub> either.

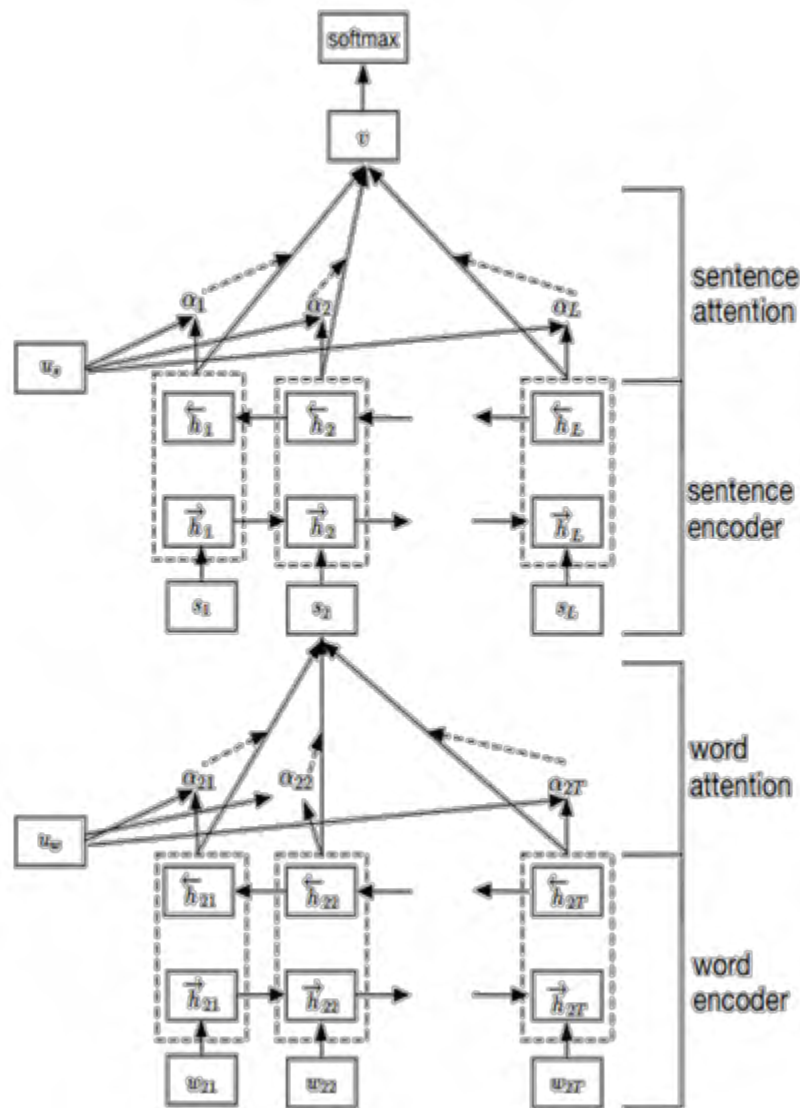
It's **not** a **very** **[good]**<sub>+3→-0.25</sub> movie song!

It **removes** my **[doubts]**<sub>-3→+1</sub>.



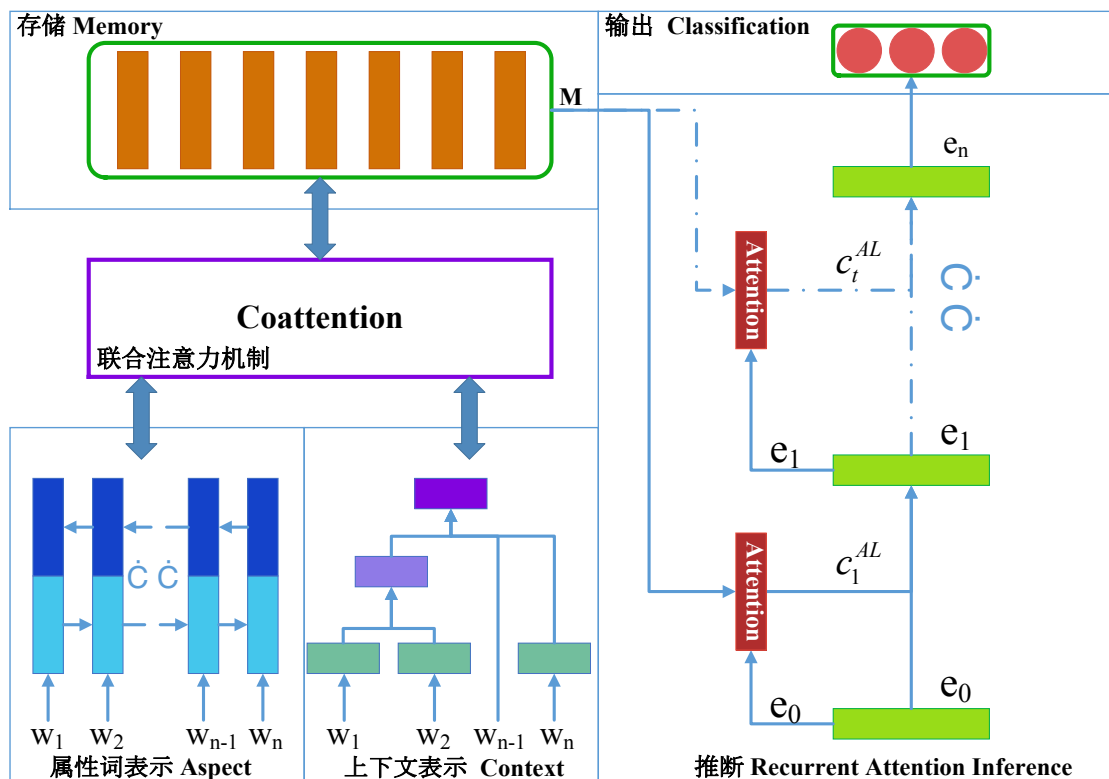
# 篇章级情感分类

- 篇章级情感分析是指对**整篇文档**的**全局情感**极性进行分析判断
- 基本思想是**对词->句子->篇章**逐层进行语义编码表示，得到篇章的向量表示
- 让机器像人一样对篇章进行阅读：从简单的字词组成句子，句子组成篇章，最后形成思想，这就是自然语言处理中的**层级 ( Hierarchical )** 概念



# 属性级情感分类

- 属性级情感分析 ( Aspect-Level ) 是指对所描述事物的属性情感极性进行判断
- 基本思想是对**属性词和其上下文**进行表示，并建立它们之间的关系，进而判断情感极性
- 基于联合注意力机制的属性级情感分析模型，对上下文和属性词分别用Tree-LSTM和Bi-LSTM，进行语义推断





# 情绪分类-细粒度的情感分析

## 情绪轮理论

- 基本情绪：8种不可以再分解的情绪
- 情绪强度：根据强弱分为3个等级
- **复合情绪**：由基本情绪叠加得来

## 复合情绪扩展

心理学家Turner的扩展复合情绪定义

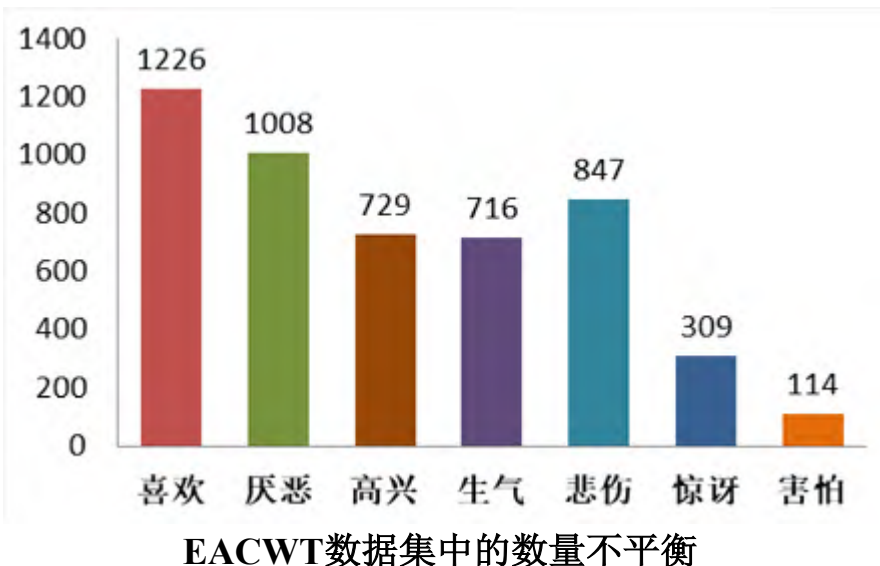
复合情绪	基本情绪成分	相对复合情绪	基本情绪成分
乐观	期待+快乐	失望	惊讶+悲伤
希望	期待+信任	难以接受	惊讶+厌恶
焦虑	期待+恐惧	愤慨	惊讶+生气
爱	快乐+信任	懊悔	悲伤+厌恶
内疚	快乐+恐惧	妒忌	悲伤+生气
惊喜	快乐+惊讶	悲观	悲伤+期待
顺从	信任+恐惧	轻蔑	厌恶+生气
好奇	信任+惊讶	愤世嫉俗	厌恶+期待
多愁善感	信任+悲伤	病态	厌恶+快乐
警惕	恐惧+惊讶	有攻击性	生气+期待
绝望	恐惧+悲伤	骄傲	生气+快乐
羞耻	恐惧+厌恶	占主导的	生气+信任



# 情绪分析存在的挑战

两个观察：

- **数量不平衡**：各个情绪类的准确率分布与样本数量大致相同
- **语义不平衡**：用户会用**相同的词**或者**类似的句式**表达不同的情绪



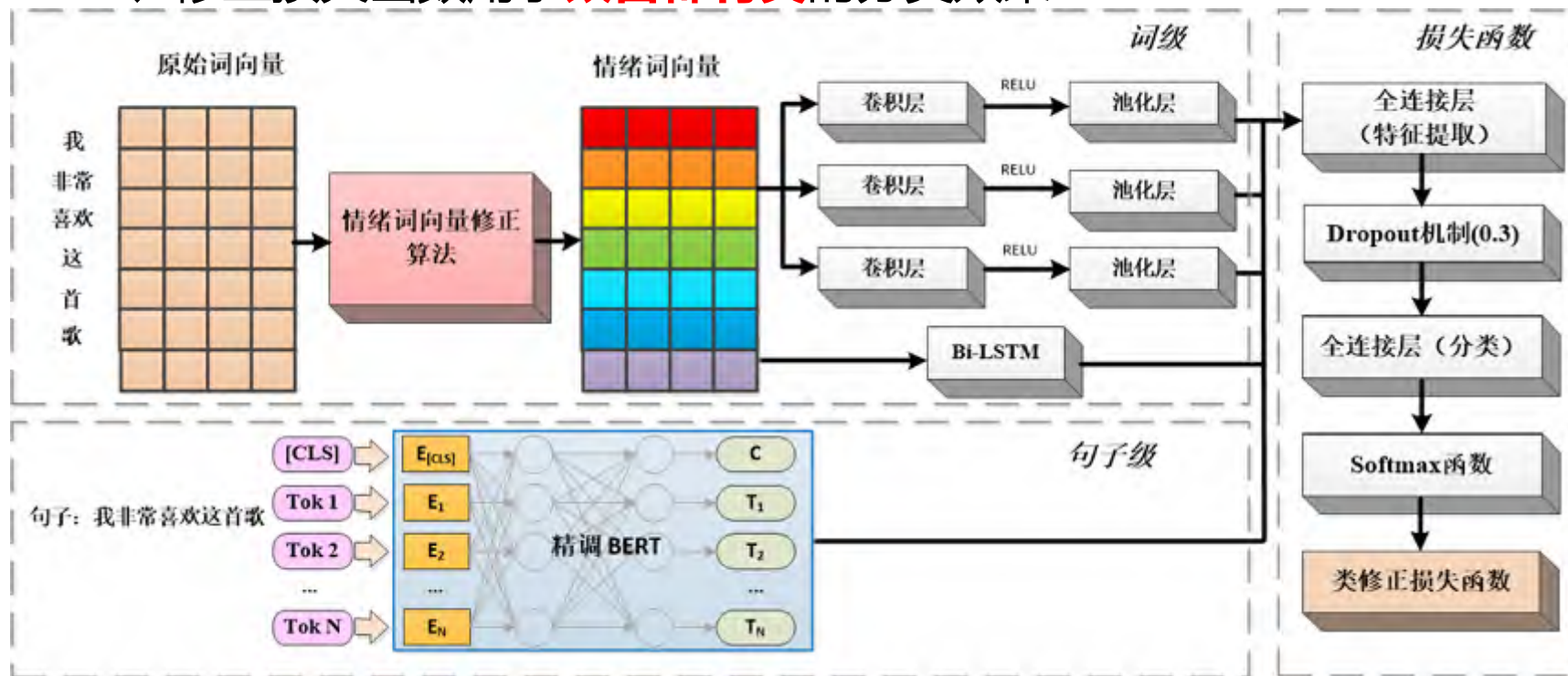
EACWT数据集中的特征不平衡

情绪	EACWT中的语料
厌恶	这两天心里特闷的慌。讨厌这种感觉。(which translates to: “These two days I feel so boring, that feel is <b>so annoying</b> .”)
生气	真是够讨厌扫兴的！(which translates to: “It's <b>so annoying</b> and disappointing.”)

# 情绪分析模型

实现框架：

- 1、词向量修正算法修正原始词向量，并使用CNN、Bi-LSTM学习**词级**情绪特征
- 2、利用情绪语料对Bert进行精调，用于学习**句子级**情绪特征
- 3、修正损失函数用于**改善稀有类**的分类效果



---

## 三、深层情感分析

# 什么是反讽识别

- 社交媒体中存在着大量的**反讽表达**，用来表达和字面意思相反的观点或情感。因此，反讽识别对于社交媒体情感分析具有重要意义。



CARAMEL\_POPCORN\_焦糖炸玉米

#香港机场离港口已全部关闭#废青加油哦废青棒棒哒

明天继续堵哦，堵到9月哦

澳门跟深圳机场能不能变成中转重要国际枢纽就看你们的了呢!



侠客岛 🇨🇳

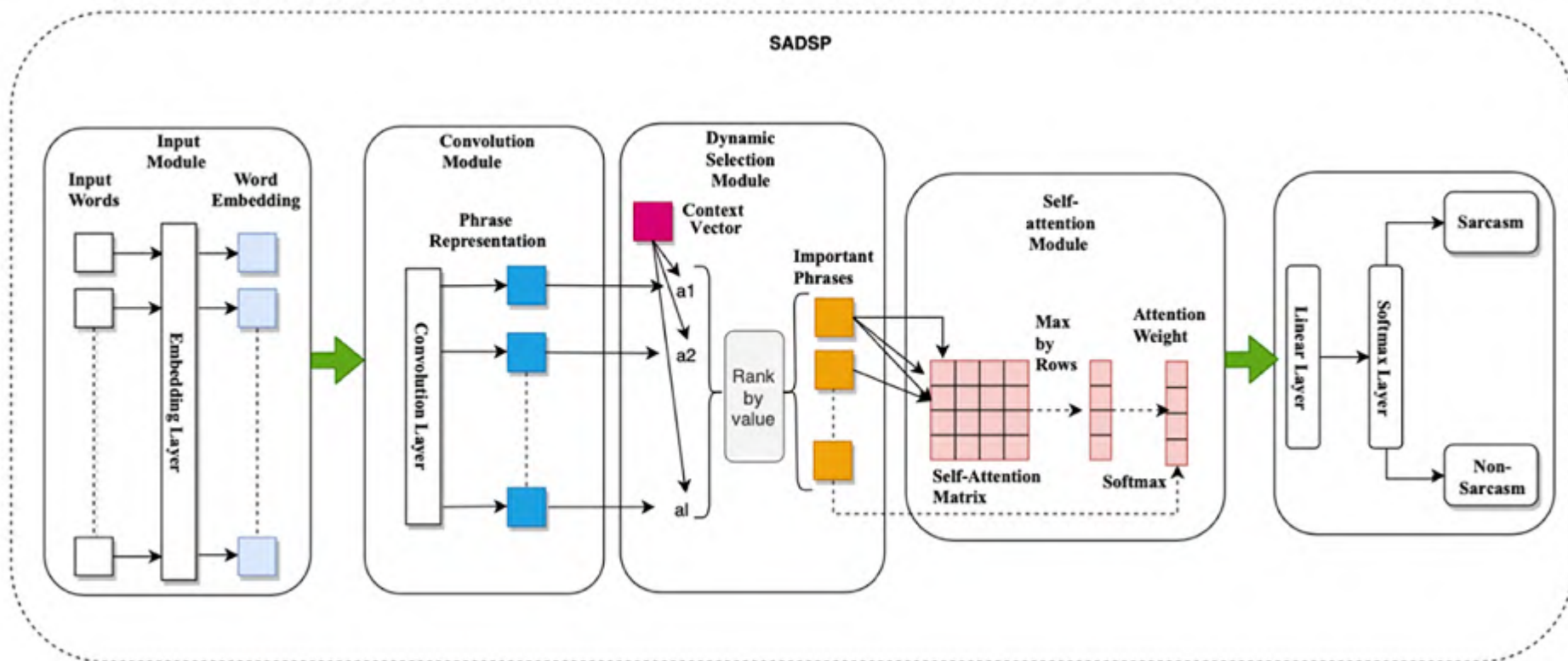
昨天 16:58 来自微博 weibo.com

【#煽动罢课的港独头子开学了#🙄】昨天，煽动“九月罢课”的港独头子罗冠聪，在社交软件上说已抵达纽约，前往耶鲁进修……在美国，他会“继承美国国务卿、国会议员”的指导，继续“展开很多工作”，并鼓动其他人继续上街。有网友反讽道，“他要去耶鲁，你要进大牢，你在为他的学位而战”。

那么问 ... 全文

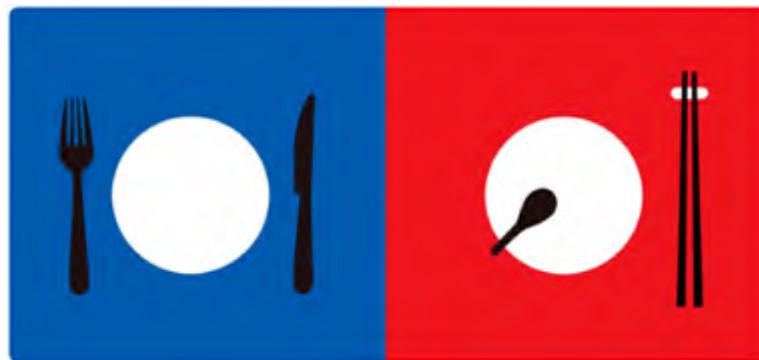
# 反讽识别模型

- 提出一个模型来捕捉句子片段间的不一致性和冲突来识别反讽，并在模型中加入动态选择机制去选择那些对于识别反讽起关键作用的片段来捕捉句子中的冲突



# 个性化情感分析

- 在情感分析中加入个性化元素（人物画像、语言、文化）



- 在情感分析中考虑社交关系、群体特性



# 多模态情感分析

单模态上的信息往往不全面或者带有歧义，多模态数据可以对单模态数据形成多视角补充。



Gorgeous!



Doing hard labor again!



"As long as the road is right, not afraid of distant future; if you believe it's worthwhile, do not care about vicissitudes."



Hello there sweetie. :)



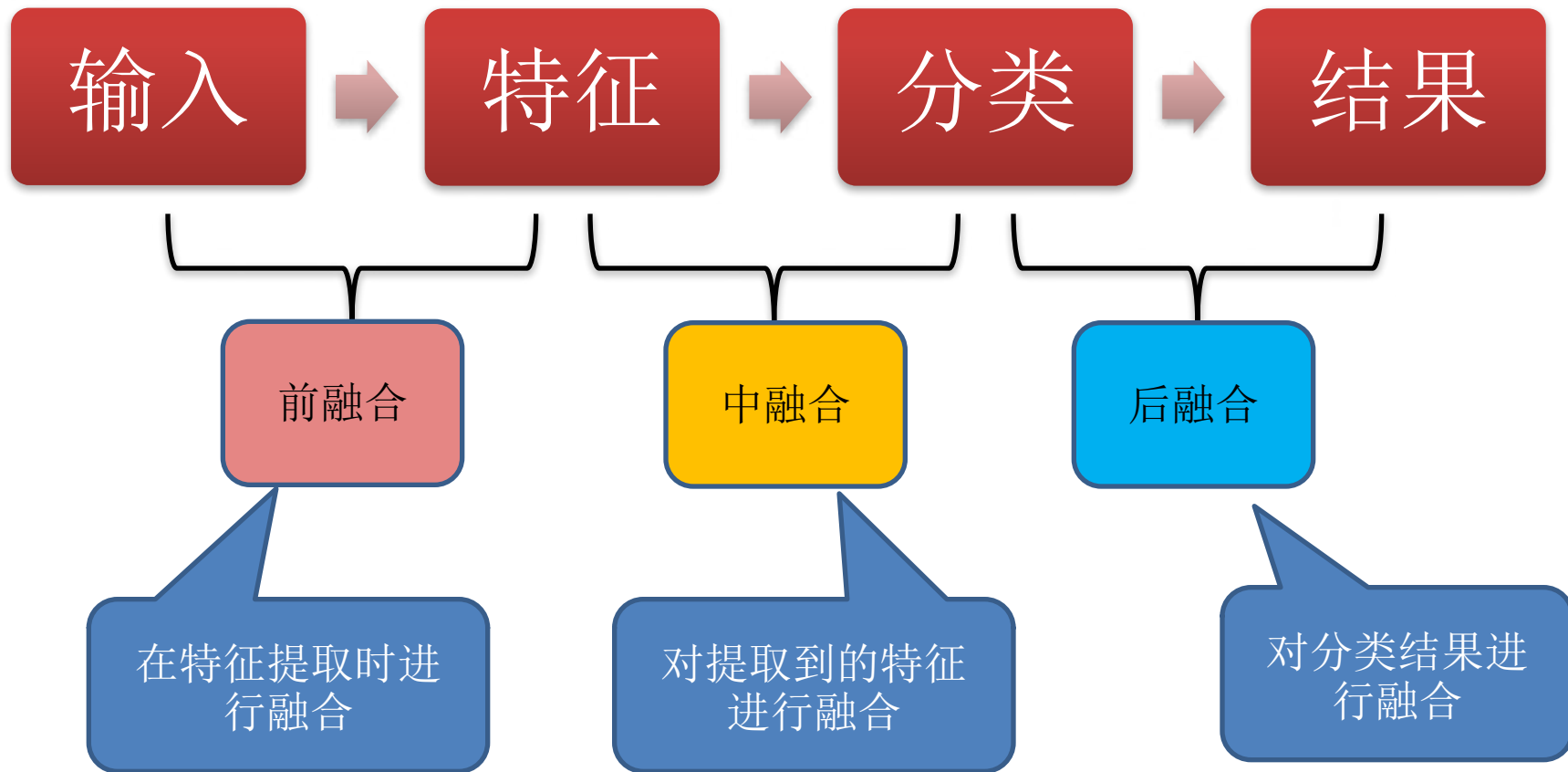
PD Achilles meets a new friend. Special post for one of our followers who I met last night and had a good chat to



If anyone woke up in edinburgh this morning to discover their car missing i think i know where it is



# 多模态情感分析



---

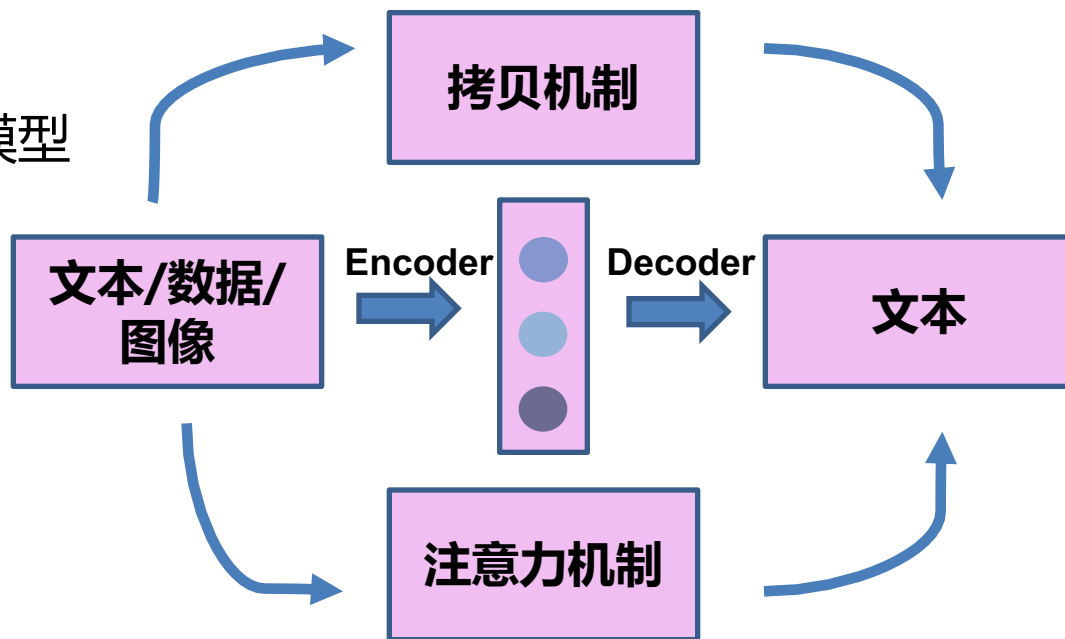
## 四、情感分析应用

# 文本生成

- 文本生成是希望计算机能够像人类一样会表达，能够撰写出高质量的自然文本
- 情感文本生成是希望机器能够生成具有情感的流畅文本

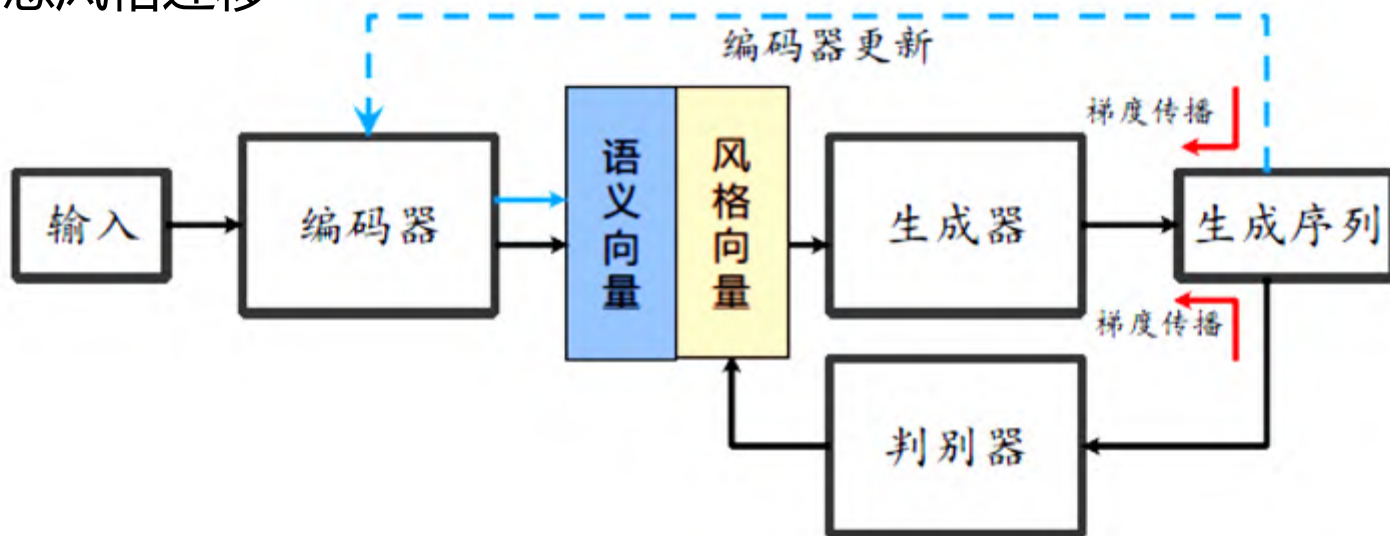
- 基于神经网络的生成模型

- 编码器-解码器
- 注意力机制
- 拷贝机制



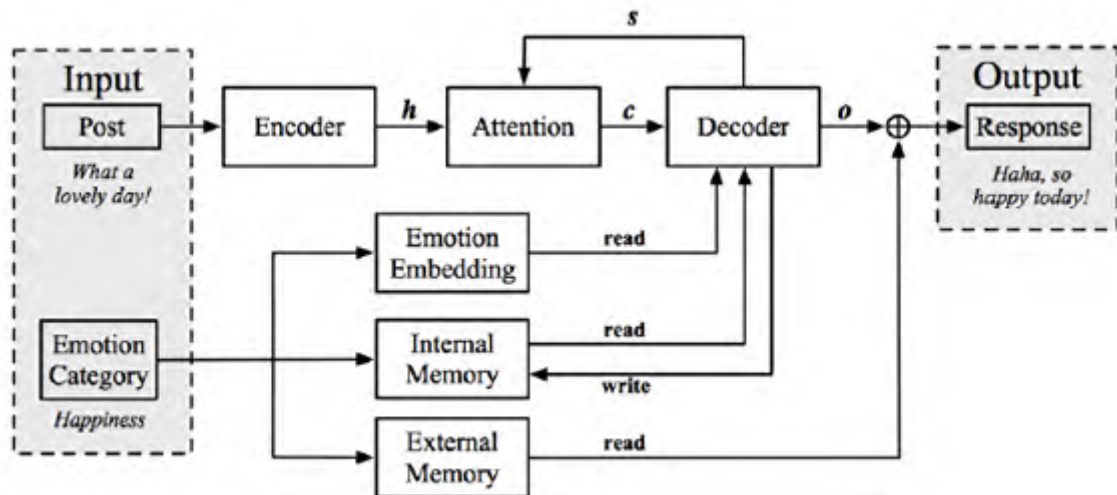
# 带情感文本生成

- 情感风格迁移



- 情绪对话生成

- 情感类别向量
- 显式情感词



# 基于情感分析的谣言检测

网络社交媒体中**充满了谣言**。而评论反映了用户所表达的情感立场。这些立场信息的表达有助于判别和辨识原始内容。例如，谣言微博与真实的新闻微博相比，在转发和评论内容中表现出更多的质疑。



时尚的晴儿1798

[新闻联播] 😊😊 人民币停产了  
🚫 央行停止印钞!!!

与时俱进，才是智者  
在这个互联网时代，秒被颠覆，  
  
你还抱着老观念，再不改变，  
将无路可走，不要老想着嘲笑别人，  
这个世界甩掉你，连声招呼都不会打[擦汗]

luzi露紫 举报了  
07月23日 08:52

小桌子哈哈 🤪 造谣  
07月08日 00:26

小瓶子c12 🤪 谣言啊  
06月29日 17:21

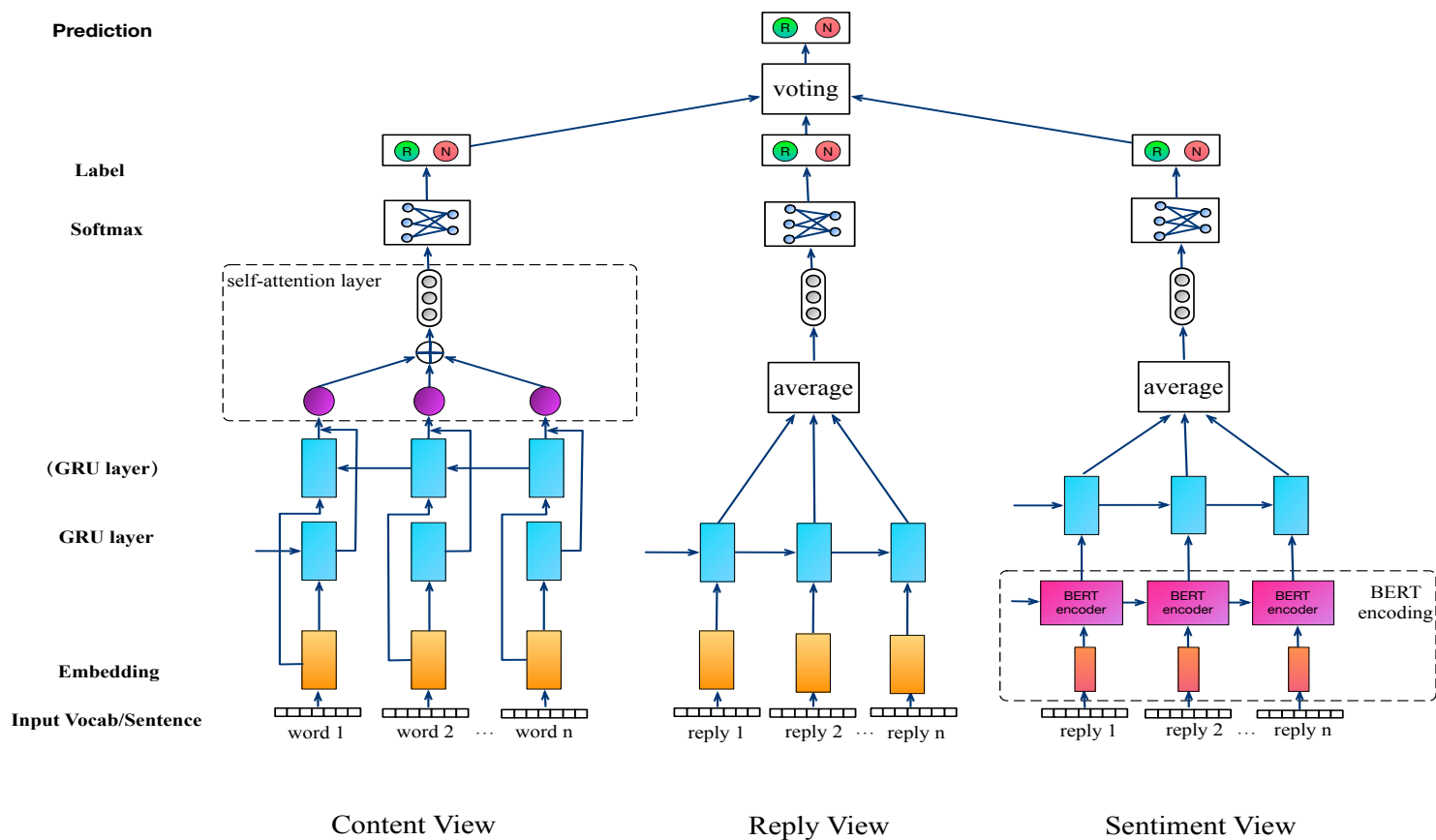
好比打篮球 🤪: @浙江网警巡查执法 造谣哦  
06月09日 23:10

野生菌六六 🤪: 回复@时尚的晴儿1798:无脑微商啧啧没救了  
06月09日 19:33

爱吃石头的麻辣小火锅 🤪: 服了，先不说是不是真的视频，丹麦央行可以宣布中国人民币停产??? 喂醒醒别睡啦  
06月09日 16:23

# 利用情感分析的谣言检测方法

- 以微博原文、评论原文和评论立场作为模型的三个不同视角，利用多个视角共同判断原文是否为谣言。



# 研究趋势

- **跨领域和跨语言情感分析**

无论是跨领域还是跨语言情感分析，都能更好利用特定语言或特定领域目前最好的模型和资源，实现情感分析在其他语言或其他领域上的迁移与适应

- **基于多媒体融合的情感分析**

结合文字、图片、语音等信息的情感分析研究将具有更好的应用前景

- **结合心理学理论的细粒度情绪分析**

情感分析大多集中于正、负二分类的粗粒度分析，结合心理学理论的细粒度情绪分析将能更好地满足实际需求

- **基于社交网络的情感分析**

结合社交关系的情感分析技术可以更好的利用群体传播影响来进行个体的情感判定

- **面向个性化的情感分析**

情感是一种高度主观的用户行为特征，为不同用户构建情感画像，可实现用户情感信息的精准化分析

# 谢谢



中国科学院 信息工程研究所  
INSTITUTE OF INFORMATION ENGINEERING, CAS